

IPMM'99

The Second International Conference on Intelligent Processing and Manufacturing of Materials

Volume 2



DISTRIBUTION STATEMENT A
Approved for Public Release
Distribution Unlimited

THIS QUALITY ASSURED
20000627 119

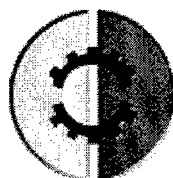
REPORT DOCUMENTATION PAGE

Form Approved
OMB NO. 0704-0188

Public Reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comment regarding this burden estimates or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302, and to the Office of Management and Budget, Paperwork Reduction Project (0704-0188,) Washington, DC 20503.

1. AGENCY USE ONLY (Leave Blank)	2. REPORT DATE June 2000	3. REPORT TYPE AND DATES COVERED Final Report
4. TITLE AND SUBTITLE IPMM'99 The Second International Conference on Intelligent Processing and Manufacturing of Materials, VOLUME 1 AND VOLUME 2	5. FUNDING NUMBERS DAAD29-99-1-0074	
6. AUTHOR(S) John A. Meech, principal investigator		
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) University of British Columbia Vancouver, BC, V6T-1Z4	8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING / MONITORING AGENCY NAME(S) AND ADDRESS(ES) U. S. Army Research Office P.O. Box 12211 Research Triangle Park, NC 27709-2211	10. SPONSORING / MONITORING AGENCY REPORT NUMBER ARO 39677.1-RT-CF	
11. SUPPLEMENTARY NOTES The views, opinions and/or findings contained in this report are those of the author(s) and should not be construed as an official Department of the Army position, policy or decision, unless so designated by other documentation.		
12 a. DISTRIBUTION / AVAILABILITY STATEMENT Approved for public release; distribution unlimited.	12 b. DISTRIBUTION CODE	
13. ABSTRACT (Maximum 200 words) The second International Conference on Intelligent Processing and Manufacturing of Materials was held in Honolulu, Hawaii on July 10--15, 1999. IPMM'99 is the second in a series of conferences dealing with the application of Artificial Intelligence and related technologies such as expert systems, fuzzy logic, artificial neural networks, genetic algorithm, pattern recognition and hybrid systems to the processing and manufacturing of materials and products. The theme of this year's conference is "Intelligence in Materials Production - The Competitive Edge!"		
14. SUBJECT TERMS		15. NUMBER OF PAGES
		16. PRICE CODE

Proceedings
of the Second International Conference
on
Intelligent Processing and Manufacturing of Materials
IPMM'99



Volume 2

Editors:
John A. Meech, Marcello M. Veiga,
Michael H. Smith, Steven R. LeClair

Hilton Hawaiian Village Hotel
Honolulu, Hawaii

July 10 - 15, 1999

DTIC QUALITY INSPECTED 4

20000627111

Proceedings of IPMM'99
The 2nd International Conference on Intelligent Processing and Manufacturing of Materials
Honolulu, Hawaii

© 1999 IEEE

IEEE Catalogue Number: 99EX296

ISBN: 0-7803-5489-3 (Softbound Edition)

Library of Congress Number: 99-61516

Copyright and Reprint Permission: Abstracting is permitted with credit to the source. Libraries are permitted to photocopy beyond the limit of U.S. copyright law for private use of patrons, those articles in this volume that carry a code at the bottom of the first page, provided the per-copy fee indicated in the code is paid through the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923. For other copying, reprint or republication permission, write to IEEE Copyrights manager, IEEE Operations Center, 445 Hoes Lane, P.O. Box 1331, Piscataway, NJ 08855-1331. All rights reserved ©1999 IEEE. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers of lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE.

printed in Vancouver, B.C. by West Coast Reproduction Centres

5000015111

IPMM'99

Foreword

It is a great pleasure to welcome you to Hawaii and to the Second International Conference on Intelligent Processing and Manufacturing of Materials.

The theme of this year's conference is

"Intelligence in Materials Production - the Competitive Edge!"

"We are living in a Material World" sings Madonna and throughout the ages, materials have been essential for bettering our standard of living. All materials derive from the Earth's crust, oceans or atmosphere and soon, even from outer space. By applying human intelligence to the properties of matter and the environment of a problem, Mankind has developed countless materials, goods and products to serve Society's needs. Perhaps Madonna's song should refer to an "Intelligent World".

IPMM'99 is the second in a series of conferences dealing with the application of Artificial Intelligence and related technologies such as expert systems, fuzzy logic, artificial neural networks, genetic algorithms, pattern recognition and hybrid systems to the processing and manufacturing of materials and products. The 1st IPMM Conference was held in 1997 in Gold Coast, Australia and attracted over 300 delegates from 37 countries with a diverse set of backgrounds that included computing, mining, metals, materials, manufacturing, etc. The participants found much to share in the "intelligent" methods being used around the world to study, simulate, process or make materials and products. The cross-disciplinary nature of this conference series is a "breath of fresh air" to many of us.

In the production of ores, minerals, metals, ceramics, plastics or food, intelligent methods have become essential to better understand and process materials or to manufacture products. Intelligence is embodied in creative ways to select components, predict properties, control processes or operate plants and factories. Such methods may be software or hardware applications; they may mimic how the human mind processes information; or they may derive from first-principle modeling of the physics and chemistry of matter.

Corporations are increasingly turning to intelligent methods to enhance their competitiveness in today's complex society and so, the technical program at IPMM'99 is focused on research aimed at leading-edge industrial applications and on the identification of newly-evolving technologies.

Intelligence exists all around us. Each of us uses it to conduct our daily lives. As the world becomes increasingly more complex and as communication systems allow massive transfer of information at incredibly reduced time scales, the global community will begin to apply this rapid collection of knowledge through powerful massively-parallel systems that currently exist within our families, communities, towns and cities, states and countries. As computers become more and more predominant in our workplaces and homes, we will begin to consider problems to which previously, we could only apply our imaginations.

Intelligence that exists in humans and other species, is now being placed into machines and materials. We are applying intelligence as we explore outer space and yes, perhaps, one day we will discover new intelligent life forms in the universe.

Conventional approaches to problem solving are becoming more and more integrated into systems that are controlled using fuzzy logic, artificial neural networks, genetic algorithms to create hybrid systems. As these systems become more widely used in industry, the complexity issues will grow as we attempt to find "optimum" solutions to our problems. You will find many papers at IPMM'99 dealing with hybrid systems that combine the attributes of many different methodologies.

The methodologies may mimic the human thought-process either symbolically or structurally. Papers are available describing evolutionary techniques that adapt to changing circumstances and allow solutions to problems to adjust in response to external factors. A number of papers focus on developing instruments that provide artificial senses that mimic the eye, the nose, the ears and yes, even the tongue. Tactile activities are also important in robotic fields and so even, the sense of touch is described in some papers.

As we examine these proceedings and its many fascinating areas of research, I wish to issue a few challenges that we face in developing new products to assist us in our future lives. Some of these ideas came to mind from reading the papers and still others developed from the difficult exercise of putting together this conference and proceedings.

Challenge 1: Can we find a way to put a film onto the surface of eye-glasses that changes its refractive index in response to external light and/or the distance at which the wearer is focusing? Perhaps, the film would have a variable R.I. from top to bottom of the lens.

Challenge 2: Can we develop hearing aids that actually work properly -- which filter out extraneous noise and provide quality hearing to those of us impaired?

Challenge 3: Can we develop a word processing program which always prints out documents the way they were originally designed regardless of the print driver and hardware being used?

The first two challenges can revolutionize the field of hearing and sight aids can improve the quality of life for many, many people. The third challenge probably exists already but is not being marketed in a way to be of widespread use. Much time and effort must be spent by those of who use word processors everyday to reformat documents as we move around our offices or as we move from home to office.

Of course, these challenges are trivial compared to some of the more fundamental (environmental, political and social) issues facing the world today. But it is even such small problems being solved that can have enormous impact on so many people. The opportunity to apply intelligence exists in everything we do or make. It is up to those of us involved in the field to see that the intelligent methods are applied for the good of Mankind.

There are many people who contributed to the success of this conference. For their advice and patience, I would like to thank the following individuals: Marcello Veiga, Mike Smith, Steve LeClair, Tom Zacharia, Guy Nicoletti, Ed Szczerbicki, Madjid Fathi, Malcolm Scoble, Tara Chandra, Lotfi Zadeh, Zoran Bugarinovic, Robert Wagoner, Junichi Endou, Susuma Shima, Debbie McCoy, Iqbal Ahmad, John Atkinson, Scott Meech, Sonia Veiga, Bojan Bugarinovic, Igor Bugarinovic. Special Mahalo (thanks) to Epoonni Perkins, Stan Omizo and Lisa Chang of the Hilton Hawaiian Village for their support and patience.

We trust you will find these proceedings of great benefit in your future endeavors and research.

John A. Meech
General Chair, IPMM'99

Honolulu, Hawaii, USA
May 31, 1999.

IPMM'99

The Second International Conference on Intelligent Processing and Manufacturing of Materials

Volume 1

Contents

Plenary Presentations	1
From Computing with Numbers to Computing with Words: From Manipulation of Measurements to Manipulation of Perceptions Lotfi A. Zadeh Berkeley Initiative in Soft Computing (BISC) Computer Science Division and the Electronics Research Laboratory, Department of Electrical Engineering and Computer Science, The University of California, Berkeley, California, USA	3
Hybrid Modeling for Testing Intelligent Software for Lunar-Mars Closed Life Support Jane T. Malin Intelligent Systems Branch, Automation, Robotics and Simulation Division NASA Johnson Space Center, Houston, Texas, USA	5
Image Analysis and Vision Systems for Processing Plants Antti J. Niemi[*], Heikki Hyötyniemi[*], and Raimo Ylinen^{**} [*] Helsinki University of Technology, Control Engineering Laboratory P.O. Box 5400, FIN-02015 HUT, Finland ^{**} University of Oulu, Systems Engineering Laboratory, P.O. Box 4300, FIN-90401 Oulu, Finland	11
Progress in Japan's Intelligent Manufacturing Systems Research Program Yuji Furukawa Tokyo Metropolitan University, Minami-Osawa, Hachioji, Tokyo, Japan	21
Analysis of Processes and Large Data Sets by a Self-Organizing Method Teuvo Kohonen Helsinki University of Technology, Neural Networks Research Centre, P.O. Box 2200, FIN-02015 HUT, Finland	27
Rough Set Theory for Intelligent Industrial Applications Zdzislaw Pawlak Institute of Theoretical and Applied Informatics, Polish Academy of Sciences, Poland	37

From Fuzzy Set Theory to Computational Intelligence – Special European Experiences	45
Hans-Juergen Zimmermann Aachen University of Technology, RWTH, Institute of Operations Research, Aachen, Germany	
Telemining™ Systems Applied to Underground Hard Rock Metal Mining at Inco Limited	53
Gregory R. Baiden INCO Mines Research, Sudbury, Ontario, Canada	
Soft Sensors for Processing Plants	59
Guillermo D. González Department of Electrical Engineering, University of Chile, Santiago, Chile	
J. Keith Brimacombe Memorial Symposium: Intelligence in Materials Engineering	71
In Memory of J. Keith Brimacombe: The Pursuit of Quality in the Casting of Materials	73
Indira V. Samarasekera The Centre for Metallurgical Process Engineering, The J.K. Brimacombe Advanced Materials and Process Engineering Laboratory (AMPEL), The University of British Columbia, Vancouver, Canada	
Towards Intelligent Steel Processing	75
Rian J. Dippenaar BHP Institute for Steel Processing and Products, The University of Wollongong, Wollongong, New South Wales, Australia	
Computer Simulation and Information Management Systems for Material Processing	85
Yoshiyuki Nagasaka Department of Distribution Science, Osaka Sangyo University, Osaka, Japan	
Simulation of Springback with the Draw/Bend Test	91
Kaiping Li, Lumin Geng, Robert H. Wagoner Dep't. Materials Science and Engineering, Ohio State University, Columbus, Ohio	
Development of an Integrated System for Designing Steelmaking Aim Compositions	105
P.A. Manohar*, S.S. Shivathaya**, M. Ferry*, T. Chandra* * Dep't. of Materials Engineering., University of Wollongong, NSW, Australia ** Hawker de Havilland Ltd., Bankstown, Australia	

A SCADA-based Expert System to Provide Delay Strategies for a Steel Billet Reheat Furnace	111
Clifford Mui*, John A. Meech**, Peter Barr**	
* Dynapro Systems Inc., Vancouver, B.C., Canada	
** The Centre for Metallurgical Process Engineering, The University of British Columbia, Vancouver, B.C., Canada	
Simulation and Analysis of Thin Strip Casting Processes	119
Yogeshwar Sahai, Manish Gupta	
Dep't. Materials Science and Engineering, Ohio State University, Columbus, Ohio	
Intelligent Manufacturing I	129
Agent-Based Control of Manufacturing Systems	131
László Monostori, B. Kádár	
Computer Automation Institute, Hungarian Academy of Sciences, Budapest, Hungary	
Intelligent Database Support for Manufacturing and Processing of Industrial Materials	139
Sylvanus A. Ehikioya, E.G. Truelove, Thomas T. Tran	
Brandon University, Brandon, Manitoba, Canada	
Intelligent Production Management in Mining Systems	145
Sean Dessureault, Malcolm Scoble, Scott Dunbar	
Dep't. of Mining and Mineral Process Engineering, University of British Columbia, Vancouver, B.C., Canada	
Intelligent Quality Control for the Food Industry using a Fuzzy-Fractal Approach	151
Oscar Castillo, Patricia Melin	
Tijuana Institute of Technology, Chula Vista, California	
Design Tool for Assessing Manufacturing Environments	157
Daniel A. Holder*, Raymond D. Harrell*, Daniel Rochoviak**, Phillip Farrington**, Dawn Russell**, John Rogers**, Sherri Messimer**	
* US Army AMCOM, Redstone Arsenal, Alabama, USA	
** University of Alabama in Huntsville, Alabama, USA	
Models, Algorithms and Decision Support Systems for Letter Mail Logistics	163
Hans-Jürgen Sebastian	
RWTH Aachen , Operations Research Group, Aachen, Germany	
Intelligent Processes for Production Control	165
Edson Pacheco Paladini	
Universidade Federal de Santa Catarina, Florianópolis, SC, Brasil	

Fuzzy Systems I	171
Industrial Applications of Fuzzy System Modeling	173
I. Burhan Turksen University of Toronto, Canada	
From Intelligent Models to Smart Ones	179
Heikki Hyötyniemi Helsinki University of Technology, Espoo, Finland	
A Fuzzy Design Evaluation Based on a Taguchi Quality Approach	185
A. Donnarumma*, N. Cappetti*, M. Pappalardo*, Esamuele Santoro ** * Università di Salerno, Italy. ** Università di Napoli, Italy	
Non-Traditional Performance Analysis	191
J. Arlen Cooper Sandia National Laboratories, Albuquerque, New Mexico, USA	
Methods of Creating Membership Functions for Fuzzy Rules in Knowledge Bases	195
Cezary Orłowski Technical University of Gdańsk, Gdańsk, Poland	
An Efficient Method for Constructing Fuzzy Rules	201
Bojan Novak University of Maribor, Stajerska, Slovenia	
Fuzzy Clustering Model Based on Changes in Vagueness	207
Mika Sato-Ilic University of Tsukuba, Ibaraki, Japan	
Thin Films and Surface Processing	213
Modeling and Control of Optical Interference Filters Using Plasma Assisted Chemical Vapor Deposition	215
Derek A. Linkens*, M.F. Abbod*, J. Metcalfe**, B. Nichols ** * University of Sheffield, Sheffield, U.K. ** GEC-Marconi Limited, Caswell, UK.	
A Study of Mechanical Properties of Multi-Layered Thin Films	221
T. Hirasawa, H. Kotera, T. Yamamoto, Y. Sakamoto, S. Shima Kyoto University, Sakyo-ku, Kyoto, Japan	
Foundations of Micro-Machining	227
Juergen Leopold Institute of Tool Engineering and Quality Management, Chemnitz, Germany	

Design of Novel Smoothing by Atomic Layer Epitaxy for Microstructure Fabrication	233
S. Hirose*, A. Yoshida**, M. Yamaura**, H. Munekata **	
*Mechanical Engineering, AIST, MITI, Tsukuba, Ibaraki, Japan	
**Tokyo Institute of Technology, Midori-ku, Yokohama, Japan	
Study of the Relationship Between Groove Cross Sectional Area per Pulse of Q-Switched Yag Laser and Strength of Processing Sound	239
T. Kurita, T. Ono	
Tokyo Metropolitan Institute of Technology, Tokyo, Japan	
Optimization of Thickness Distribution of Micro-Membrane by Genetic Algorithm	245
Hidetoshi Kotera, Y. Sakamoto, T. Hirasawa and S. Shima	
Kyoto University, Sakyo-ku, Kyoto, Japan	
Manufacturing of Metallic Prototypes and Tools by Laser Cutting and Diffusion Bonding	251
S. Sändig, P. Wiesner	
Dep't. of Mechanical Engineering, Technical University of Ilmenau, Germany	
Evolutionary Systems and Machine Learning	255
Artificial Immune Systems: a New Frontier in Artificial Intelligence	257
Dipankar Dasgupta*, Stephanie Forrest **	
* University of Memphis, Tennessee, USA	
** University of New Mexico, Albuquerque, NM, USA	
Inductive Learning for Optimization of Simulation Model Output	269
Rainer Barton*, Helena Szczerbicka **	
* German Aerospace Center (DLR), Institute for Flight Mechanics, Braunschweig, Germany	
** University of Bremen, Bremen, Germany	
A Genetically Optimised Fuzzy Parser of Natural Language	277
Olgierd Unold	
Wroclaw University of Technology, Wroclaw, Poland	
A Genetic Algorithm-based Approach to Solve Process Plan Selection Problems	281
K.M. Tiwari*, S.K. Tiwari*, Debjit Roy*, N.K. Vidyarthi **.	
* Manufacturing Engineering, National Institute of Foundry and Forge Technology, Hatia, Ranchi, India	
** Mechanical Engineering, NERIST, Nirjuli, Itanagar, India	
*** SriVenkateshNagar, Chennai-600092, India	

Breeding Policies in Evolutional Approximation of Optimal Subspace H.M. Huang and P.L. Leung City University of Hong Kong, Kowloon, Hong Kong	285
Prediction of Cement Paste Mechanical Behaviour from Chemical Composition using Genetic Algorithms and Artificial Neural Networks José C. Cassa, Giovanni Floridia, André R. Souza, Rodrigo T. Oliveira Universidade Federal da Bahia, Salvador, Bahia, Brazil	291
Rough Sets-based Machine Learning Using a Binary Discernability Matrix Reynaldo Felix, Toshimitsu Ushio Systems and Human Science, Osaka University, Toyonaka, Japan	299
Intelligence in the Design of Materials and Processes I	307
INTELLIGOLD - An Expert System For Gold Plant Process Design Vanessa Torres*, Arthur Torres**, John A. Meech *** *Companhia Vale do Rio Doce, Belo Horizonte, Brazil **University of Sao Paulo, SP, Brazil ***University of British Columbia, Vancouver, Canada	309
A Hardware Design for Real-Time Multiple Target Tracking Frederick Ferguson, Chandra Curtis North Carolina A&T State University, Greensboro, NC, USA	317
Low-Cost Supersonic Missile Inlet Fabrication Technique C.S. Cornelius, D.A. Gibson US Army Aviation and Missile Command, Redstone Arsenal, AL	325
Design of High Performance Missile Structures Utilizing Advanced Composite Material Technologies J.R. Esslinger, R.N. Evans, G.W. Snyder US Army Aviation and Missile Command, Redstone Arsenal, AL	331
Modelling the Mechanical Stability of Metal Catalyst Carriers C. Guist, H. Bode Bergische Universität-Gesamthochschule Wuppertal, Germany	339
Integration of Newly Developed AI Assembly, Production, and Material Flow Virtual Tools Daniel A. Holder*, Raymond D. Harrell*, Terri L. Calton**, John F. Atkinson*, Brandy M. Brasfield* * US Army AMCOM, Redstone Arsenal, Alabama ** Sandia National Laboratories, Albuquerque, NM	347

Prediction of Materials Properties	353
How ab-initio Computer Simulation Can Predict Materials Properties Before Experiment Yoshiyuki Kawazoe Tohoku University, Sendai, Japan	355
Data Driven Knowledge Extraction of Materials Properties J.S. Kandola*, S.R. Gunn*, I. Sinclair**, P.A.S. Reed** University of Southampton, U.K.	361
A Quantum Neural Net: with Applications to Materials Science B. Igel'nik*, M. Tabib-Azar*, Y.-H. Pao*, and S. R. LeClair** *Case Western Reserve University, Cleveland, OH, USA **Material Directorate, Wright Laboratory, Fairborn, OH, USA	367
Ontology for Phase Diagram Databases N. Ono, R. Kainuma, H. Ohtani, K. Ishida, M. Kato Tohoku University, Sendai, Japan	373
Prediction of Concrete Mechanical Behaviour from Data at Lower Ages using Artificial Neural Networks José C. Cassa, Giovanni Floridia, André R. Souza, Rodrigo T. Oliveira Universidade Federal da Bahia, Salvador, Bahia, Brazil	381
Improving the Prediction Accuracy of a Constitutive Model with ANN Models L.X. Kong and P.D. Hodgson Deakin University, Geelong, Victoria, Australia.	389
Hybrid Fuzzy Modelling Using Simulated Annealing: Application to Materials Property Prediction Min-You Chen, Derek A. Linkens The University of Sheffield, Sheffield, UK	395
Intelligence in Materials Science I	401
Inorganic Glasses: Old and New Structures on the Eve of the 21st Century J. Šesták*, B. Hlaváček⁺, N. Koga*** * Czech Academy of Sciences, Prague, Czech Republic ** University of Pardubice, Pardubice, Czech Republic *** Hiroshima University, Higashi-Hiroshima, Japan	403
Oxygen Solubility Modeling in Aqueous Solutions Desmond Tromans University British Columbia, Vancouver, B.C., Canada	411

- On the Oxidation of Steel in CO₂ and Air** 417
Gity Samadi Hosseinali, Ainul Akhtar
 Powertech Labs Inc., Surrey, British Columbia, Canada
- Retardation of Hydrogen Embrittlement by Electrolytic ZrO₂ Coating of AISI 430 Stainless Steel** 423
I.B. Huang, S.K. Yen
 National Huwei Institute of Technology, Huwei, Taiwan
- The Effect of Ca Addition on Viscosity and Electrochemical Properties of Mg-Alloys Produced by Casting** 429
H.S. Kim*, Shuji Hanada*, Ha-Guk Jeong*, Dong-Wha Kum **
 * Tohoku University, Sendai, Japan
 ** Korea Institute of Science and Technology, Seoul, Korea
- Bio-Compatible Ceramics as Mimetic Material for Bone Tissue Substitution** 431
Zdenek Strnad*, Jaroslav Šesták **
 *Lab. for Glass and Ceramics (LASAK), Prague, Czech Republic
 **Czech Academy of Sciences, Prague, Czech Republic
- Intelligent Design of GaSb doped Single Crystals** 437
B. Štěpánek, J.Šesták, J.J.Mareš, J.Křištofik, V.Šestáková, P.Hubík
 Academy of Sciences of the Czech Republic, Semiconductor Department, Prague, Czech Republic,
- Intelligent Image Analysis Applications** 443
- Astronomical Image Processing - Applications To Ultra-Faint Imaging of Small, Moving, Solar System Bodies: Comets and Near-Earth Objects** 445
Karen J. Meech
 University of Hawaii, Institute for Astronomy, Honolulu, HI, USA
- A High Performance Computing Algorithm for Improving In-Line Holography** 447
Hesham Eldeib
 Electronic Research Inst., National Research Center, Giza, Egypt
- Human Face Detection System by KenzanNET with Preprocess Analyzing of Hyperspectral Image** 453
Takakazu Chashikawa*, Keizo Fujii, Yoshiyasu Takefuji ***
 * Keio University, Kanagawa, Japan.
 **NITTAN Co., Ltd., Japan

Using Image Analysis and Partial Least Squares Method to Estimate Mineral Concentrations in Mineral Flotation Jari Hätönen* , Heikki Hyötyniemi* , J. Miettunen** , L.-E. Carlsson*** *Helsinki University of Technology, Espoo, Finland **Outokumpu Mining Oy, Pyhäsalmi Mine, Pyhäsalmi, Finland ***Boliden Mineral AB, Mineral Processing, Boliden, Sweden	459
A Combined Morphological and Color-Based Approach to Characterize Flotation Froth Bubbles Giuseppe Bonifazi, Silvia Serranti, F. Volpe, R. Zuco Università degli Studi di Roma, "La Sapienza", Italia	465
Robust Bubble Delineation Algorithm for Froth Images Weixing X. Wang, O. Stephansson Department of Civil and Environmental Engineering, Royal Institute of Technology, Stockholm, Sweden	471
The Characterization of Flotation by Colour Information and Selecting the Proper Equipment A.K. Sirén VTT Information Technology, Espoo, Otaniemi, Finland	477
Intelligence in Environmental Applications	479
Robust Engineering Approaches to Maximize Results in Business, Cost, Engineering, Human, Quality and System Technologies Roberto C. Villas Bôas ** CYTED - Science and Technology for Development in Iberoamerica, Mineral Technology Sub-Program, Madrid, Spain	481
Imaging Techniques for Process Optimization and Control in Glass Recycling Giuseppe Bonifazi, Paolo Massacci Ingegneria Chimica, dei Materiali, Materie Prime e Metallurgia, Università degli Studi di Roma "La Sapienza", Roma, Italia	485
Application of Heuristic Modeling in Natural Resource Sciences Steven Mackinson Fisheries Centre, University of British Columbia, Vancouver, B.C., Canada	491
ARDEX - A Fuzzy Expert System for ARD Site Remediation Judita Balcita, John A. Meech University of British Columbia, Vancouver, B.C., Canada	499

Modeling of Gold Heap Leaching for Criteria of Sustainability Targets	505
Luiz R. P. De Andrade Lima*, Roberto C. Villas-Bôas**	
* Federal University of Bahia, Salvador, BA, Brazil	
** Center for Mineral Technology, Rio de Janeiro, RJ, Brazil	
Design Optimisation of Aluminium Recycling Using the Taguchi Approach	513
A.R. Khoei, D.T. Gethin, I. Masters	
Mechanical Engineering, University of Wales Swansea, UK	
Towards a Better Understanding of Environmental Science through Application of Fuzzy Sets	519
Mory M. Ghomshei, John A. Meech	
University of British Columbia, Vancouver, B.C., Canada	
Intelligence in Rolling Processes	527
Data Mining and State Monitoring in Hot Rolling	531
<u>L. Cser</u>* **, A.S. Korhonen**, P.Mäntylä***, O. Simula**, J.Ahola **	
* Bay Zoltan Institute for Logistics and Production Technology, Miskolc-Tapolca, Hungary	
** Helsinki University of Technology, Espoo, Finland	
*** Rautaruukki Steel, Raase, Finland	
Determination of Thickness Control Parameters of Rolling Processes by the Sensitivity Method, using Neural Networks	537
<u>Luis E. Zárate</u>*, Horacio Helman **	
* Departamento de Ciência da Computação, Pontifícia Universidade Católica de Minas Gerais, Belo Horizonte, Brazil	
** Departamento de Engenharia Metalúrgica e de Materiais, Universidade Federal de Minas Gerais, Belo Horizonte, Brazil	
AI Approach to Modeling Rolling Loads in Design of Cold Rolling Processes	543
J. Kusiak*, J.G. Lenard**, K. Dudek*	
* Akademia Gorniczo-Hutnicza, Krakow, Poland	
** University of Waterloo, Waterloo, Ontario, Canada	
Direct Determination of Sequences of Passes for Strip Rolling Process by Means of Fuzzy Logic Rules	549
C.D.M.Pataro, H. Helman	
Universidade Federal de Minas Gerais, Belo Horizonte, Brazil	
Elongation-Controlled Rolling of H-Shaped Wire	555
H. Utsunomiya, M. Shinkawa, F. Shimaya, Y. Saito	
Materials Science and Engineering, Osaka University, Japan.	

Application of a Neural Network to Speed Up a Mathematical Model for Calculation of Strip Profiles in Flat Rolling Yukio Shigaki, Horacio Helman Universidade Federal de Minas Gerais, Belo Horizonte, Brazil	561
Intelligent Methods in Metal Forming Processes	563
A Fundamental Study of the Incremental Deep Drawing Process S. Shima, H. Kotera, K. Kamitani, S. Nagatomo Mechanical Engineering, Kyoto University, Kyoto, Japan	565
Intelligent Design Architecture for Process Control of Deep-Drawing K. Manabe*, H.Koyama*, K.Kato**, S. Yoshihara *** * Mechanical Engineering, Tokyo Metropolitan University, Japan ** Integrated Systems Japan,Ltd., Tokyo, Japan *** Tokyo National College of Technology, Tokyo, Japan	571
An Iterative Approach to Determine Heat-Treatment and Composition from the Mechanical Yield Strength of an Al-Li Alloy James M. Fragomeni, Ohio University, Mechanical Engineering, Athens, Ohio, USA	577
A Design of Experiments Statistical Approach to Determine the Effect of Extrusion Process Variables on the Mechanical Properties of a Heat-Treated Al-Li Alloy James M. Fragomeni Ohio University, Mechanical Engineering, Athens, Ohio, USA	585
Control of Liquid Segregation of Semi-Solid Al-Alloys during Intelligent Compression Testing C.G. Kang, K.D. Jung, H.K. Jung Pusan National University, Mechanical Engineering, Korea	593
Adaptability to Frictional Change of Fuzzy Adaptive Blank Holder Control for Deep Drawing S. Yoshihara*, K. Manabe**, H. Nishimura ** * Tokyo National College of Technology, Mechanical Eng., Japan ** Tokyo Metropolitan University, Tokyo, Japan	601
An AI Process Control System with Simulation Database and Adaptive Filter for V-Bending M. Yang*, A. Katayama*, K. Manabe*, N. Aikawa ** * Dep't. of Mechanical Engineering, Tokyo Metropolitan University, Japan ** Tokyo Engineering University, Japan	607

Intelligent Manufacturing II	613
The Distributed Intelligent Control of Complex Systems	615
Wayne J. Davis University of Illinois at Urbana-Champaign, Department of General Engineering, Urbana, IL, USA	
PDM-based Virtual Enterprises – Bridging the Semantic Gap	623
A. Karcher, J. Wirtz Dep't. of Mechanical Engineering, Technical University of Munich, Garching, Germany	
A Methodology to Diagnose the Target Cost in a Manufacturing Process	629
A. Ariotti, C. Fantozzi, M. Granchi, E. Vettori Mechanical, Nuclear and Manufacturing Engineering, University of Pisa, Italy	
Resource Allocation for a Fast-Tracked Project	635
Yassiah Bissiri, Scott Dunbar Department of Mining and Mineral Process Engineering, University of British Columbia, Vancouver, B.C., Canada	
Hybrid Simulation Objects using Fuzzy Set Theory for Simulation of Innovative Process Chains	641
T. Menzel, M. Geiger Dep't. of Manufacturing Technology, University Erlangen-Nuremberg, Erlangen, Germany	
Manufacturing Management Improvement through Rapid Production of Budgets	649
E.J. Colville School of Engineering, University of Tasmania, Hobart, Tasmania, Australia	
A Connectionist Method to Solve Job Shop Problems	655
Marko Fabiunke, Gerd Kock GMD Research IT Center,(FIRST) Berlin, Germany	
Fuzzy Systems II	661
Designing in Many-Valued Logic	663
A. Donnarumma and Michele Pappalardo University of Salerno, Mechanical Engineering, Fisciano, Italy	

Modulus Genetic Algorithm and its Application to Fuzzy System Optimization Sinn-Cheng Lin Dep't. of Educational Media and Library Sciences, Tamkang University, Tamsui, Taipei Hsien Taiwan, PRC	669
Fuzzy Evolutionary Programming for Portfolio Selection in Investment Decisions T. Van Le Faculty of Information Sciences and Engineering, University of Canberra, Belconnen, Australia.	675
Design of a Region-Wise Fuzzy Sliding Mode Controller with Fuzzy Tuner C.C. Kung, W.C. Lai Dep't. of Electrical Engineering, Tatung Institute of Technology, Taipei, Taiwan	681
A Multi-Input Current-Mode Fuzzy Integrated Circuit for Pattern Recognition Gu Lin, Bingxue Shi Institute of Microelectronics, Tsinghua University, Beijing, PRC	687
A Framework for Intelligent Systems based on Vector-Annotated Logic Programs Kazumi Nakamatsu*, Yumi Hasegawa*, Jair Minoru Abe**, Atsuyuki Suzuki ***. *Himeji Institute of Technology, School of Humanity, Environment Policy and Technology, Himeji, Hyogo, Japan **Paulista University, Sao Paulo, Brazil. ***Shizuoka University, Japan	695
A Fuzzy Logic Assisted Electrodynamic Balance for Unit Operations on Single Levitated Particles M. Pappalardo*, A. Pellegrino*, M. d'Amore**, P. Giordano**, P. Russo ** * Dep't. of Mechanical Engineering, University of Salerno, Fisciano, Italy ** Dep't. of Chemical and Food Engineering., University of Salerno, Fisciano, Italy	703
Author's Index	I-1

IPMM'99

The Second International Conference on Intelligent Processing and Manufacturing of Materials

Volume 2

Contents

Artificial Neural Networks I	711
Artificial Neural Networks (ANN) as Simulators and Emulators: An Analytical Overview	
Guy M. Nicoletti University of Pittsburgh at Greensburg, Pennsylvania, USA	713
Logical Rule Extraction from Data by Maximum Neural Networks	
T. Saito, Y. Takefuji Keio University, Fujisawa, Kanagawa, Japan	723
Iterative RBF Neural Networks as Metamodels of Stochastic Simulations	
George Meghabghab, George Nasr Dep't. of Mathematics and Computer Science, Valdosta State University, Valdosta, Georgia, USA	729
A Systematic and Reliable Approach to Pattern Classification	
R.Doraiswami, M.Stevenson, S. Rajan Department of Electrical Engineering, University of New Brunswick, Fredericton, New Brunswick, Canada	735
Dynamic Associative Memory Using Chaotic Neural Networks	
Yoshihisa Fukuhara, Yoshiyasu Takefuji Keio University, Graduate School of Media and Governance, Fujisawa, Kanagawa, Japan	743
Trends in Intelligent Process Control in the Primary Aluminium Industry	
R.T. Bui*, L. Tikasz*, J. Perron ** * Université du Québec a Chicoutimi, Chicoutimi, Québec, Canada ** Alcan International Ltd, Jonquiere, Quebec, Canada	749
Modeling of the Flow Stress Relationship using a BP Network	
Y. Y. Yang, Derek A. Linkens University of Sheffield, Sheffield, U.K.	755

Intelligence in Materials Science II	763
The Heredity and Control of Microstructures of Liquid Metals During Rapid Cooling Processes	765
Rang-su Liu, Ji-yong Li, Hai-rong Liu	
*Department of Physics, Hunan University, Changsha, P.R. China	
**Department of Chemistry, Hunan University, P.R. China	
AI Approach to Internal Variable-based Rheological Model for Steels	773
J. Kusiak, M. Pietrzyk	
Department of Metallurgy and Materials Engineering, Akademia Gorniczo-Hutnicza, Krakow, Poland	
Electrolytic ZrO₂ Coating on Co-Cr-Mo Implant Alloys of Hip Prosthesis	779
S.K. Yen, H.Z. Zan, M.J. Guo	
National Chung Hsing University, Taichung, Taiwan	
Automated Stress Control of Electroplated Nickel-Phosphorus Alloy	785
G.T. Yu, M. Williams	
National Huwei Institute of Technology, Materials Eng., Taiwan	
Mechanism of Electrolytic Coating of Al₂O₃ on MAR-M247 Superalloy	789
S.K. Yen, C.C. Chang	
Dep't. of Material Engineering, National Chung Hsing University, Taichung, Taiwan	
A New Process to Produce Advanced Zirconia-based Ceramic Composites from Low-Value Minerals	797
Sonia M. B.Veiga*, Marcello M. Veiga**, A.C.D. Chaklader**, J. C. Bressiani *	
* Inst. de Pesquisas Energéticas e Nucleares , São Paulo, Brasil.	
** University of British Columbia, Vancouver, BC, Canada	
High Temperature Flow Stress Model and Hot Deformation Behaviors for High-Mo Austenitic Stainless Steel	805
Xu Yourong, Chen Liangshen, Jin Lei, Wang Deying	
Materials Science and Engineering, Shanghai University, Jiading, P.R. China.	
Intelligent Manufacturing III	811
Information Management of Complex Systems: Perspectives for the New Millennium	813
Z. Gomolka*, E. Szczerbicki **	
* University of Szczecin, Szczecin, Poland	
** University of Newcastle, Newcastle, Australia	

Present Status of Intelligent Machines in Sheet Metal Fabricating and Forming in Japan Junichi Endou Kanagawa Institute of Technology, Kanagawa, Japan	817
Design of Enterprise Network Communication Subsystems Adam Grzech Wroclaw University of Technology, Poland	823
The Industrial Desktop – Real Time Business and Process Analysis to Increase Productivity in Industrial Plants Osvaldo A. Bascur OSI Software, Inc., The Woodlands, Texas, USA	829
Enterprise Staff Scheduling by Genetic Algorithm Search Tiehua Zhang*, William A. Gruver*, Michael H. Smith ** * Simon Fraser University, Burnaby, BC, Canada **University of California, Berkeley, CA, USA and	839
Intelligence in Surface Processing of Materials	845
Intelligent AE Sensor for Monitoring of Finish Machining Process Slavko Dolinšek*, J. Kopac*, Z.J. Viharos*, L. Monostori** * University of Ljubljana, Mechanical Engineering, Slovenia ** Hungarian Academy of Science, Budapest, Hungary	847
A New Fuzzy-Fractal Approach for Surface Quality Control in Intelligent Manufacturing Of Materials P. Melin, O. Castillo Tijuana Institute of Technology, Chula Vista, California, USA	855
A Study on Axisymmetric Indentation by the Rigid-Plastic Finite-Boundary Element Method Yong-Ming Guo, Kenji Nakanishi Dep't. of Mechanical Engineering, Kagoshima University, Kagoshima, Japan	861
Design of Intelligent Spindle for High Speed Machining B.L. Zhang, Y.P. Li, B.S. Zhu, P. Ma, Y. Luo Guangdong University of Technology, Guangzhou, China	867
Robotics and Intelligent Control I	869
Autonomous Control of Complex Dynamical Systems in Support of a Manned Mission to Mars James A. Kurien, Daniel J. Clancy NASA Ames MS 269-3, Moffett Field, California, USA	871

Mining Automation in the Next Millennium: a Tele-Operated LHD Vehicle Model	879
Yeen-Shien Hwang^{*1}, Neda Farmer^{**2}, Jason Hart **	
[*] University of British Columbia, Vancouver, B.C., Canada	
¹ Huckleberry Mines Ltd., Houston, B.C., Canada	
² Luscar Coal Mine, Hinton, A.B., Canada	
** Nautilus International Limited, Burnaby, B.C., Canada	
Dynamic Reconfiguration of Holonic Lower Level Control	887
X. Zhang*, D.H. Norrie*, A. Kusiak **	
[*] University of Calgary, Canada	
** University of Iowa, USA	
Intelligent Process Monitoring for Paper Machines	895
Janos L. Grantner*, Peter E. Parker**, George A. Fodor ***	
[*] Department of Electrical and Computer Engineering, West Michigan University, Kalamazoo, Michigan, USA.	
** Department of Paper and Printing Science and Engineering, West Michigan University, Kalamazoo, Michigan, USA.	
*** ABB Automation Products AB, Vasteras, Sweden	
An Integration Design Approach in PID Controller	901
Jen-Yang Chen	
China Institute of Technology and Commerce, Taipei, Taiwan	
Holonically Object Oriented System	909
Shigeki Sugiyama	
Gifu Industry and Technology Research Center, Kasamatsu-Cho, Hashima-Gun, Gifu-Ken, Japan.	
Intelligent Instrumentation and Measurement	919
Pre-Processing of Industrial Process Data for Outlier Detection and Correction	921
Jonathan Tenner*, Derek A. Linkens*, T.J. Bailey **	
[*] University of Sheffield, UK.	
**British Steel Engineering Steels U.K. Ltd.	
Intelligent Measurement System Confirmation	927
P. H. Osanna, M.N. Durakbasa	
Vienna University of Technology, Wien, Austria	
Simulation of the Dynamic Properties of Nuclear Meters in Coal Preparation Control Systems	933
Stanislaw Cierpisz	
Silesian Technical University, Poland	

Acoustic Emission Monitoring of SAG Mill Performance	939
S.J. Spencer, J.J. Campbell, K.R. Weller, Y. Liu CSIRO Minerals, Queensland, Australia	
Novel Polymeric Electrochemical/Chemical Sensors and Display Devices Integrated with Artificial Intelligence	947
A. Talaie***, J.Y.Lee***, Y.K. Lee****, J. Jang*, D.J. Choo****, S.H. Park****, G. Huh****, J.A. Romagnoli** * Physics Department, Kyung Hee University, Seoul, Korea ** Chemical Engineering, Sydney University, Sydney, Australia *** Chemistry Dept., NSW University, Sydney, Australia **** Kyung Hee University, Seoul, Korea	
Material Properties under Drawing and Extrusion with Cyclic Torsion	953
L.X. Kong, P.D. Hodgson, L. Lin and B. Wang School of Engineering and Technology, Deakin University, Geelong, Victoria, Australia	
Artificial Neural Networks II	959
Neural Network-based Resistance Spot Welding Quality Prediction	961
N. Ivezic, J.D. Allen, Jr., T. Zacharia Oak Ridge National Laboratory, Oak Ridge, TN, USA	
An Adaptive Artificial Neural Network to Model a Cu/Pb/Zn Flotation Circuit	967
Saiedeh Forouzi, John A. Meech University of British Columbia, Vancouver, B.C., Canada	
Multivariable Predictive Neuronal Control Applied to Grinding Plants	975
Manuel Duarte*, Alejandro Suárez**, Danilo Bassi *** * Dep't. Ing. Eléctrica, Universidad de Chile, Santiago, Chile ** Dep't. de Electrónica, Univ. T.F. Sta. María, Valparaíso, Chile *** Dep't. de Informática, Univ. de Santiago, Santiago, Chile	
Practical Neural Network Applications in the Mining Industry	983
Logan Miller-Tait, Rimas Pakalnis University of British Columbia, Vancouver, B.C., Canada	
Neural Network Model and Model-Based Control of Deformation Processing	989
Nenad Ivezic, John D. Allen, Jr., Thomas Zacharia Oak Ridge National Laboratory, Oak Ridge, TN, USA	

Verifying Detected Facial Parts by Multidirectional Associative Memory	995
Miki Kitabata, Yoshiyasu Takefuji Keio University, Fujisawa Kanagawa, Japan	
A Current-Mode Sorting Circuit for Pattern Recognition	1003
Gu Lin , Bingxue Shi Microelectronics, Tsinghua University, Beijing, P.R. China	
Intelligence in the Design of Materials and Processes II	1009
Intelligent Design Methods for Smart Materials	1011
Madjid Fathi-Torbaghan, L. Hildebrand Dep't. of Computer Science, University of Dortmund, Dortmund, Germany	
Identification of a Model Which Relates Variations in Shape Geometry to Process Control Variables of Shape Forging	1017
B.F. Rolfe*, M.J. Cardew-Hall*, G.A.W. West**, S.M. Adballah* *Australian National University, Canberra, ACT, Australia **Curtin University of Technology, Perth, Australia	
Mechanical Characteristics of HIPed SiC Particulate-Reinforced Al-Alloy MMCs	1023
C.Y. Chung, K.C. Lau City University of Hong Kong, Kowloon, Hong Kong, P.R. China	
Hydrostatic Extrusion of Composite Rod	1029
Ui-Bin Tsai, Chi-Wei Wu, Ray-Quen Hsu National Chiao-Tung University, Hsin-Chu, Taiwan	
Numerical Modelling and Localized Failure Analysis in Metal Powder Forming Processes	1035
A.R. Khoei, R.W. Lewis, D.T. Gethin Mechanical Engineering, University of Wales Swansea, UK	
Microstructure and High Temperature Deformation Behavior of a TiN/Ti₅Si₃Nano-Grain Composite Produced by Non-Equilibrium PM Processing	1041
Kei Ameyama and Yasuhiko Suehiro Ritsumeikan University, Kusatsu City, Shiga, Japan	
Shape Prediction of Growing Billet in Spray Casting using a Scanning Gas Atomizer	1047
Eon-Sik Lee*, Sangho Ahn* and Shinill Kang ** *Research Institute of Industrial Science and Technology, Advanced Materials Division, Pohang, Kyungbuk, South Korea **Yonsei University, Seoul, Korea	

Intelligence in Concurrent Engineering	1053
Modelling Design Planning in Concurrent Engineering C. Reidsema and E. Szczerbicki Dep't. of Mechanical Engineering, University of Newcastle, NSW, Australia	1055
Computer-Aided Integrated Design for Injection Molding Yuh-Min Chen*, Rong-Shean Lee*, ChengTer Ted Ho *** * National Cheng Kung University, Tainan, Taiwan *** National KaoHsiung Inst. of Science and Technology, Taiwan	1061
Artificial Psychology – an Attainable Scientific Research on the Human Brain Zhiliang Wang, Lun Xie University of Science & Technology (USTB), Beijing, P.R. China	1067
Soft-Object Technology for Autonomous Manufacturing Components Control Ahmed Hambaba College of Engineering, San Jose State University, San Jose, CA	1073
A Monitoring Framework for Software Project Development Ho-Leung Tsoi* and Derek Cheung ** * Software Quality Institute, Griffith University, Australia ** Computer Studies, City University of Hong Kong, Hong Kong	1079
Redefining the Web: toward the Creation of Large-Scale Distributed Applications Guy M. Nicoletti Engineering Department, University of Pittsburgh at Greensburg, Pennsylvania, USA	1087
How Can We Form/Expand Conceptions in Workers' Minds According to Their Individualities? Kumiko Ishino Konan University, Utsunomiya-City, Kobe, Japan	1093
Robotics and Intelligent Control II	1101
Navigation by Weighted Chance S. Reimann*, A. Mansour ** *German National Research Center for Information Technology, Birlinghoven, Germany ** Bio-Mimetic Control Research Center, (RIKEN), Nagoya, Japan	1103

Vehicle Routing Problem Using Clustering Algorithm by Maximum Neural Networks	1109
Noriko Yoshiike, Yoshiyasu Takefuji Keio University, Fujisawa, Kanagawa, Japan	
Acquisition of Communication Protocol for Autonomous Multi-AGVs Driving	1115
Michiko Watanabe, Masashi Furukawa Asahikawa National College of Technology, Hokkaido, Japan	
Heuristic Neuro-Fuzzy Model For Evaluation of Urban Transportation Projects	1123
Marcus Vinicius Quintella Cury, Saul Fuks Universidade Federal do Rio de Janeiro, UFRJ, Brasil	
Optimal Controller Design for Finite Word Length Implementation using a Genetic Learning Algorithm	1125
Wen-Shyong Yu Tatung Institute of Technology, Taipei, Taiwan	
Adaptive Fuzzy Controller for Non-Linear Uncertain Systems	1131
Chiang-Cheng Chiang, Chih-Chien Hu Tatung Institute of Technology, Taipei, Taiwan, R.O.C.	
Hybrid Modeling (view-graphs)	1137
Holistic Strategies for Designing Multistage Material Processes	1139
W.G. Frazier Air Force Research Laboratories, Wright-Patterson AFB, Ohio	
A New Methodology of Using Design of Experiments as a Precursor to Neural Networks for Material Processing: Extrusion Die Design	1151
<u>Bhavin Mehta</u>, Hamza Ghulman, Rick Gerth Ohio University, Athens, Ohio	
Incorporating Hybrid Models into a Framework for Design of Multi-Stage Material Processes	1157
E. Medina Air Force Research Laboratories, Ohio	
Hybrid Modeling for the Interdisciplinary Design of More Affordable Systems	1163
<u>J. Poindexter</u>, Gerald R. Shumaker, Brian A. Stucke Air Force Research Laboratories, Ohio	

Hybrid Modeling for Testing Intelligent Software for Lunar-Mars Closed Life Support Jane T. Malin Intelligent Systems Branch, Automation, Robotics and Simulation Division NASA Johnson Space Center, Houston, Texas, USA	1179
Discrete Modeling via Function Approximation Methods - Towards Bridging Atomic- and Micro-Scales <u>A.G. Jackson</u>, M. Benedict Air Force Research Laboratories, Dayton, Ohio, USA	1185
Microstructure Predictions From Atomistic Rule Set Cellular Automata M.O. Zacate, <u>R.W. Grimes</u>, P.D. Lee Imperial College, London University, London, England, U.K.	1197
Fuzzy Molecular Modeling David A. Ress North Carolina State University, North Carolina, USA	1225
Imaging Studies and Density Functional Analysis of Surfaces and Interfaces: Comparison of Theory and Experiment <u>John F. Maguire</u>, Steven R. LeClair Air Force Research Lab, Wright-Patterson AFB, Dayton, Ohio	1235
Modeling Gas Byproducts from MOCVD Thin-Film Depositions <u>J. G. Jones</u>, P.D. Jero Air Force Research Lab, Wright-Patterson AFB, Dayton, Ohio	1241
Imaging for Process Optimization and Control (view-graphs)	1247
Nondestructive Imaging of Surface & Sub-Surface Defects in Thin-Films with Super Spatial Resolution using Evanescent Microwave Fields Massood Tabib-Azar Case Western Reserve University, Cleveland, Ohio	1249
Investigation of Raman Imaging for Advanced Control of YBCO Cool- Down Processing using Pulsed Laser Deposition <u>J.D. Busbee</u> *¹, R.R. Biggers*, J.G. Jones*, D.V. Dempsey*², G. Kozlowski ** * AFRL, Materials, Wright-Patterson AFB, Dayton, Ohio ** AFRL, PRP, Wright-Patterson AFB, Ohio ¹ Technical Management Concepts, Inc., Beavercreek, Ohio ² University of Dayton Research Institute, Dayton, Ohio	1258
Process Control Via Gaze Detection Technology <u>Jaihie Kim</u>*, Gang Ryung Park*, Steven R. LeClair ** * Yonsei University, Seoul, Korea ** AFRL, Wright-Patterson AFB, Dayton, Ohio	1263

- The Third Eye Cameras - Dynamic and Static Hyperspectrum Imaging** 1271
Yoshiyasu Takefuji
 Environmental Information, Keio University, Fujisawa, Japan
- The Third Eye Approach to Innovative Designs and Applications into the 21st Century - Human Recognition System by Nonlinear Oscillations** 1277
Souichi Oka, Yoshiyasu Takefuji, William Huang
 Environmental Information, Keio University, Fujisawa, Japan
- Intelligent Rate Control for MPEG-4 Coders** 1285
Gwanh Hoon Park, Jae Hyung Park, Yoon Jin Lee
 Yonsei University, Computer Science, Wonju, Kwangwon, Korea
- Concept, Development, Mass Production, and Applications of Artificial Retina Chips** 1297
Kazuo Kyuma
 Mitsubishi Electric Corporation, Japan
- Data Reduction via Auto-Associative Neural Networks** 1305
Claudia Kropas-Hughes,
 Air Force Research Lab, Wright-Patterson AF Base, Dayton, OH
- Image Processing Plume Fluence for Superconducting Thin-Film Depositions** 1317
J.G. Jones*, R.R. Biggers*, J.D. Busbee¹, D.V. Dempsey², G. Kozlowski **
 * Air Force Research Lab, Materials Direct., WPAFB, Dayton, OH
 ** Air Force Research Laboratory, PRP, WPAFB, Dayton, OH
¹ Technical Management Concepts, Inc. Beaver Creek, OH
² University of Dayton Research Institute Dayton, OH
- Innovations in Materials Design** 1321
- Towards the Future: Innovations in Materials Design** 1323
Suichi Iwata
 RACE, University of Tokyo, Faculty of Eng., Hongo, Japan
- Atomic Environments in Relation to Compound Prediction** 1339
Jo Daams*, Pierre Villars **
 * Phillips Research, The Netherlands
 ** Materials Phases Data System (MPDS), Vitznau, Switzerland
- Analysis and Visualization of Category Membership Distribution in Multivariate Data** 1361
Yoh-Han Pao*, B.F. Duan*, Y.L. Zhao*, Steven R. LeClair *
 * Case Western Reserve University, Cleveland, Ohio
 ** Air Force Research Lab, Wright-Patterson AFB, Dayton, OH

Whitney Reduction Networks for Process Discovery	1371
Mark Oxley Mathematics and Statistics, Air Force Inst. Tech., WPAFB, Ohio	
Algorithms for Predicting Properties of Materials from Intelligent Materials Design by Hyperspace Data Mining	1381
Nianyi Chen*, <u>Dongping Daniel Zhu</u> ** * Shanghai Metallurgy Inst. of Chinese Acad. of Sci., P.R. China ** Zaptron Systems, Inc., Mountain View, CA, USA	
Data Bases and Semantic Networks for Inorganic Materials Computer Design	1387
N.N. Kiselyova A.A.Baikov Institute of Metallurgy and Materials Science, Russian Academy of Science, Moscow, Russia	
First-Principles Calculations for Materials Science: Their Power and Limitations	1397
Wanda Andreoni IBM Research Division, Zurich Research Laboratory, Switzerland	
Interplay Between Large Materials Databases, Semi-Empirical Approaches, Neuro-Computing and First Principle Calculations	1399
<u>Pierre Villars</u> *, Steven R. LeClair **, Suichi Iwata *** * Material Phases Data System (MPDS), Vitznau, Switzerland ** Air Force Research Laboratory, Wright-Patterson AFB, Ohio *** RACE, Faculty of Eng., University of Tokyo, Hongo, Japan	
Software Package "MATERIALS DESIGNER" and its Application in Materials Research	1417
<u>Nianyi Chen</u>*, Wencong Lu**, Ruiliang Chen*, Pei Qin * Shanghai Metallurgy Inst. of Chinese Acad. of Sci., P.R. China ** Department of Chemistry, Shanghai University, P.R. China	
Author's Index	II-1

Artificial Neural Networks I

Artificial Neural Networks (ANN) as Simulators and Emulators - An Analytical Overview

Guy M. Nicoletti

Engineering Department, University of Pittsburgh at Greensburg
1150 Mt. Pleasant Road, Greensburg, Pennsylvania

ABSTRACT

Because of their ability to exploit the tolerance for imprecision and uncertainty in real-world problems, and their robustness and parallelism, artificial neural networks (ANNs) and their techniques have become increasingly important for modeling and optimization in many areas of science and engineering. As a consequence, the market is flooded with new, increasingly technical software and hardware products. This paper presents an analytical overview of the most popular ANNs, both in hardware and software modes. The paper is organized as follows. **Part I** introduces a basic overview of ANNs. **Part II** discusses global optimization for ANN training, the NOVEL hybrid method is presented and its performance is discussed. **Part III** discusses the techniques and means for parallelizing neurosimulations of ANNs, both at a high programming level and at a low hardware-emulation level. **Part IV** presents vector microprocessor architectures and the Spert II fixed-point system as applied to multimedia and human-machine interface. Finally, **Part V** introduces the most recently explored concept of cellular neural networks (CNN), its performance and operation are analyzed. Conclusions and recommendations conclude the paper.

INTRODUCTION

Modern digital, von Neumann type, parallel computers outperform humans in the domain of numerical computations and related symbol manipulation. The course of long evolutionary process has provided the human brain with many superior characteristics not present in von Neumann architectures. These include: *massive parallelism, distributed representation and computation, learning ability, adaptivity, inherent contextual information processing, fault tolerance, and low energy consumption.* Inspired by biological neuronets, ANNs are massively parallel computing systems consisting of a very large number of simple processors with many interconnections. Thus, ANN models attempt to simulate some "organizational" principles present in the human brain. The structure and architecture of biological neuron and neural nets are well formulated in the related literature. To place the state of ANNs in perspective, it is worth mentioning that the cerebral cortex contains about 10^{11} neurons. Each neuron is connected to 10^3 to 10^4 other neurons and, in total, the human brain contains approximately 10^{14} to 10^{15} interconnections. Messages are modulated on pulse train variable frequency of a few to several hundred Hertz, which is a million times slower than switching speeds of modern electronic circuits. However, complex perceptual decisions are typically made by a network of neurons within a few hundred milliseconds. Thus, the human brain runs parallel programs of about 100 steps long (*the 100 step rule*) [1].

Computational Model of a Neuron. The mathematical neuron as proposed by McCulloch and Pitts [2], computes a weighted sum of n input signals, x_j , $j = 1, 2, \dots, n$, and generates an output of 1 if this sum is above a certain threshold u . Mathematically,

$$y = \theta \left[\sum_{j=1}^n w_j x_j - u \right] \quad 1.$$

Where θ is a unit step function at 0, and w_j is a synapse weight associated with the j th input. This model, however, contains a number of simplifying assumptions that do not reflect the true behavior of biological neurons. Expression (1) is commonly generalized by using an activation function other than the threshold function such as the sigmoid function defined by:

$$g(x) = 1 / (1 + \exp \{-\beta x\}) \quad 2.$$

where β is the slope parameter.

Network Architectures. Based on the connection pattern (architecture), ANNs are grouped into two categories:

- **feed-forward** networks, in which graphs have no loops. Generally speaking, these networks are *static*, i.e., they produce only one set of output values rather than a sequence of values from a given input. Their response to an input is independent of the previous network state.
- **recurrent** (or feedback) networks, in which loops occur because of feedback connections. These networks are *dynamic* systems. As a new input pattern is presented, the neuron outputs are computed. With the feedback paths, the inputs to each neuron are then modified, which leads the network to enter a new state.

Different network architectures require appropriate learning algorithms.

Learning. Learning ability is the fundamental trait of intelligence. The learning process in ANNs consists of the problem of updating network architecture and connection weights so that the network can efficiently perform a specific task. Usually, the network must learn the connection weights from available training patterns. Performance is improved over time by iteratively updating the weights into the network. ANNs ability to automatically *learn from examples* constitutes one of the major advantages of such networks over traditional expert systems. The design of a learning process hinges on two fundamentals: a *learning paradigm*, i.e. the type of information available to the network, and the type of *learning rules* that govern the updating process. The procedure for implementing the learning rules is referred to as a *learning algorithm*. There are three main learning paradigms: *supervised*, *unsupervised*, and *hybrid*. There are four basic types of learning rules: error-correction, Boltzmann, Hebbian, and competitive learning.

Error-Correction Rule. These are based on the error signal $(d - y)$ used to modify the connection weights to gradually reduce the error. Here, y is the actual output generated by the network, and d is the desired output. The perceptron, (a single neuron ANN with adjustable weights), learning rule is based on this principle. Its learning algorithm is as follows:

1. *Initialize the weights and threshold to small random numbers.*
2. *Present a pattern vector $(x_1, x_2, \dots, x_n)^T$ and evaluate the output on.*
3. *Update the weights according to:*

$$w(t + 1) = w_j(t) + \eta(d - y) x_j \quad 3.$$

where d is the desired output, t is the iteration number, and η ($0.0 < \eta < 1.0$) is the gain (step size).

Boltzmann Learning Rule. This type of learning can be viewed as a special case of error-correction learning. The error is measured as the difference between the correlations among the outputs of two neurons under constrained and free-running operating conditions. Boltzman machines are symmetric ($w_{ij} = w_{ji}$), recurrent networks consisting of binary units (+1 for "on" and -1 for "off"). The rule dictates that the change in connection weight w_{ij} is given by

$$\Delta w_{ij} = \eta (\rho'_{ij} - \rho_{ij}) \quad 4.$$

where h is the learning rate, and ρ'_{ij} and ρ_{ij} are the connections between the states of units i and j when the network operates in constrained mode and free-running mode respectively. The values of ρ'_{ij} and ρ_{ij} are usually estimated from Monte Carlo methods, which is very slow.

Hebbian Rule. Based on Hebb's *postulate of learning* [3]. This postulate is derived from observations in neurobiological experiments: If neurons on both sides of a synapse are activated synchronously and repeatedly,

the synapse strength is selectively increased. Mathematically we obtain:

$$w_{ij}(t+1) = w_{ij}(t) + \eta y_j(t) x_i(t), \quad 5.$$

where x_i and y_j are output values of neurons i and j , respectively, which are connected by the synapse w_{ij} , and η is the learning rate.

Competitive Learning Rules. Unlike Hebbian learning, competitive-learning output units compete among themselves for activation. Thus, only one output unit is active at any given time. This phenomenon, (found to exist in biological neural networks), is known as *winner-take-all*. The simplest competitive learning network consists of a single layer of output units. Each output unit i in the network connects to all the input units (x_j) via weights w_{ij} , $j = 1, 2, \dots, n$. Each output unit also connects to all other output units via inhibitory weights but has a self-feedback with an excitatory weight. As a result of competition, only the unit i^* with the largest (or the smallest) net input becomes the winner, i.e., $w_{i^*} \cdot x \geq w_i \cdot x, \forall i$, or $\|w_{i^*} - x\| \leq \|w_i - x\|, \forall i$. When all the weight vectors are normalized, the two inequalities are equivalent. Mathematically, the competitive rule can be stated as:

$$\Delta w_{ij} = \begin{cases} \eta(x_j^* - w_{ij}), & i = i^* \\ 0, & i \neq i^* \end{cases} \quad 6.$$

It can be seen from (6) that the network doesn't stop learning (updating weights) until the learning rate η is 0.

Multi layer Feed-Forward Networks Generally, a standard L -layer feed-forward network (excluding the layer of input nodes) consists of an input stage, $(L-1)$ hidden layers, and an output layer of units sequentially connected (fully or locally) in a feed-forward fashion with no connections between units in the same layer and no feedback connections between layers.

Multi layer Perceptron. In this network, [4], each computational unit employs either the *thresholding* function or the *sigmoid* function. The learning algorithm generally used is the back-propagation algorithm: a gradient descend method to minimize the squared-error cost function in the equation:

$$E = \frac{1}{2} \sum_{i=1}^p \|y^{(i)} - d^{(i)}\|^2 \quad 7.$$

Radial Basis Function Network. Radial Basis Function (RBF) networks, [5], is a special case of Multi layer feed-forward networks. It has only two layers. Each unit in the hidden layer employs a radial basis function, such as a Gaussian kernel, as the activation function. The radial basis function is centered at the point specified by the weight vector associated with the unit. There are a variety of learning algorithms for the RBF network. *The basic one, which converges faster than the back propagation algorithm, employs a two-step learning strategy, or hybrid learning* [4]. However, for many problems, the RBF often involves a large number of hidden units, its speed is often slower, and as for the Multi layer perceptron, it is problem dependent, and has the same asymptotic approximation power as the Multi layer perceptron.

There are many issues in designing feed-forward networks. These include: (1) the number of layers needed for a given task, (2) the number of units needed per layer, (3) generalization ability, i.e. the network performance with data not included in the training set and, (4) the optimal size of the training set for "good" generalization.

Kohonen's Self-Organizing Maps. The self-organizing map (SOM), [6], has the desirable property of topological preservation: the important aspect of the feature maps in the cortex of highly developed animal

brains. It basically consists of a two dimensional array of units, each connected to all n input nodes. If \mathbf{w}_{ij} denotes the n - dimensional vector associated with the unit at location (i,j) of the 2D array, then each neuron computes the Euclidean distance between the input vector \mathbf{x} and the stored weight vector \mathbf{w}_{ij} .

Kohonen's SOM can be used for projection of multi-variate data, density approximation, and clustering. It has been successfully applied in the areas of image processing, robotics, and process control [7]. The design parameters include the *dimensionality of the neuron array*, the *number of neurons in each dimension*, the *shape of the neighborhood*, the *shrinking schedule of the neighborhood*, and the *learning rate*.

Adaptive Resonance Theory Model. Adaptive Resonance Theory Models (ART1, ART2, and ARTMap) were developed by Carpenter and Grossberg to create networks capable of overcoming the *stability-plasticity* dilemma: *how to learn new things (plasticity) and yet retain stability to ensure that existing knowledge is not erased or corrupted*. The complexities of such models precludes their presentations in this paper. The reader is referred to related literature on ART1, and ART2 models [8], [2].

Hopfield Network. The Hopfield network uses a network energy function as a tool for designing *recurrent networks* and for understanding their dynamic behavior [9]. Its formulation makes explicit the principle of storing information as dynamically stable attractors in the use of recurrent networks. Applications especially directed to associative memory and to the solution of combinatorial optimization problems. A Hopfield network with n units has two versions: binary and continuously valued. The network dynamics for the networks are:

$$v_i = \text{Sgn} \left[\sum_j w_{ij} v_j - \theta_j \right] \quad 8$$

Where, v_i , the output of the i th unit, is either +1 or -1, but for continuous networks, it can be any value between 0 and 1. w_{ij} is the synapse weight on the connection from unit i to unit j .

The energy function of the binary network in a state $\mathbf{v} = (v_1, v_2, \dots, v_n)^T$ is given by:

$$E = \frac{1}{2} \sum_i \sum_j w_{ij} v_i v_j \quad 9.$$

The central property of the energy function is that a network state evolves according to the network dynamics (equation 8), the network energy always decreases and eventually reaches a local minimum point (attractor) where the network remains with a constant energy.

GLOBAL OPTIMIZATION

The algorithms described in the previous section find their roots in function-minimization algorithms that can be classified as local-or global-minimization algorithms. They focus on either extreme-local search or global search. Such algorithms, in general, do not work well. A recent proposal, [10], formulates a *hybrid* method, called *Nonlinear Optimization via External Lead* (NOVEL), that combines local and global searches to explore the solution space, locate promising regions, and find local minima. In exploring the *solution space*, the search is guided by a continuous terrain-independent *trace* that does not get trapped in local minima. In locating *promising regions*, NOVEL uses a local gradient to attract the search to a local minimum. Finally, one initial point is selected for each promising local region. These points are then used for a descent algorithm to find local minima. Good global-minimization methods are: *simulated annealing* (SA), *evolutionary algorithms* (Eas), *cascade correlation with multi starts* (Cascor-MS), *gradient descent with multi starts* (Grad-MS), *truncated Newton's method with multi starts* (TN-MS). Benchmark problems studied so far, [11], are: *Two-spiral problem* – to discriminate between two sets of training points that lie on two distinct spirals in the x-y plane; *Sonar problem* – to discriminate between sonar signals bounced off a metal and a rock; *10- parity problem* – to train a network that computes the binary sum of 10 binary digits; and the *NetTalk* problem – to train a network to produce proper *phonemes*, given a string of letters as input. Results indicates that NOVEL represents a significant advance in supervised learning of feed-forward neural networks and optimization of general high-dimensional nonlinear continuous functions.

SIMULATING ANNs ON PARALLEL ARCHITECTURES

ANNs can be implemented as a *simulation* programmed on a general-purpose computer or as an *emulation* realized on special purpose hardware. Although simulations offer comfortable software environments for developing and analyzing ANNs, the computational needs of realistic applications exceed the capabilities of sequential computers. Parallelization became therefore, necessary to cope with the high computational and communication demands of *neuro-applications* [11]. Neurosimulations can be parallelized in several ways.

The amount of parallelism achieved depends on the granularity of problem decomposition. The most popular techniques are illustrated in Fig. 1.

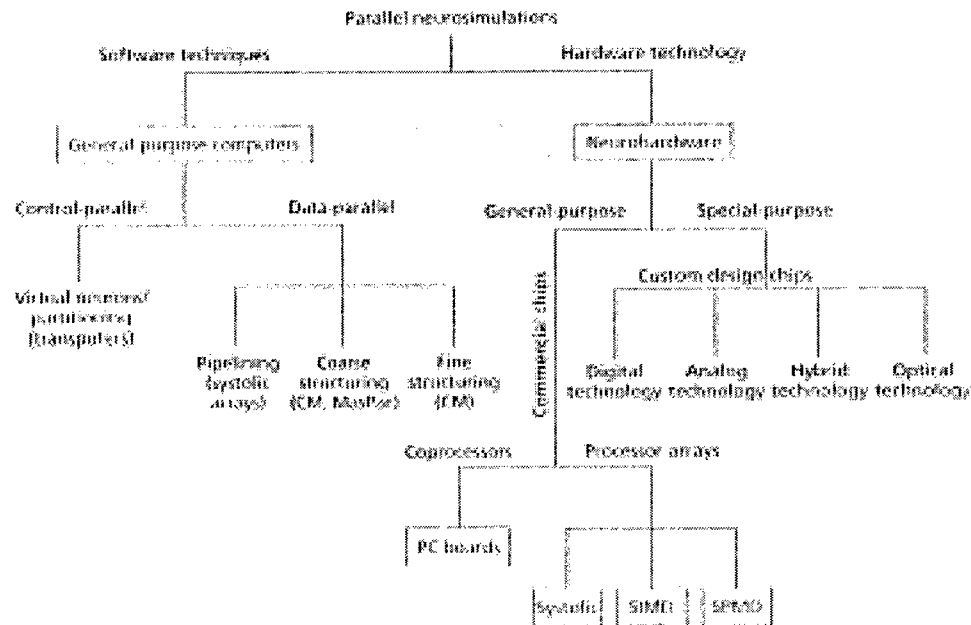


Fig. 1. A taxonomy of parallelization approaches for neurosimulations.

The parallel techniques implemented on general-purpose computers are presented in Table 1.

Table 1. Parallel implementations on general-purpose computers.

Structuring technique	No. of processors	Computer architecture	Benchmark	Performance*	
				CPS	CUPS
Coarse: training, node per layer	64K	Connection Machine ³	NETtalk	180M	38M
Coarse: training, node per layer	16K	MasPar ⁴	NETtalk	176M	42M
Fine: node, weight	64K	Connection Machine ²	NETtalk		13M
Pipelined: training, layer	10	Warp ⁵	NETtalk		17M
Pipelined: layer, node	13K	Systolic array ⁶	NETtalk		248M
Coarse: partitions	6	Transputers ⁷			207K

*The results are from the late eighties and early nineties. The performance figures are not impressive today, but the implementation techniques are still being used.

Finally, Table 2. presents recent commercially available neurocomputers.

Table 2. Commercially available neurocomputers.

Product name	Architecture/ technology	Software technique	Capacity		Performance ^a	
			Neurons	Connections	CLIPS	CPS
Synapse	Systolic	Pipelining		None	33M	5.12G
CNAPS	SIMD	Node per layer	64	128K	1.46G	5.70G
RAP	SPMD	Node per layer			106M	574M
SAIC	PC + board	ANSim; ANSpec	Virtual	Virtual	2M	11M
Balboa	PC + board	ExploreNet			9M	25M
Encoro	Transputer + board		64	512	4.2M	19M
WSI	Digital		576	36K	2.30G	
ETANN	Analog		64	10K	None	2G
Boltzmann	Hybrid		336	28K	28G	1T

^aPerformance figures should be seen as relative, since the presented systems used different precision and were not tested with a common benchmark.

General purpose neurohardware appears to offer an optimal solution, achieving both efficiency and flexibility at an acceptable price. Special purpose neurohardware demonstrates the best performance (see IEEE Micro, June 1994). The best performance is achieved with special-purpose neurocomputers that implement a particular neural model directly in silicon. Current R & D efforts are directed toward neurocomputers that consist of several modular components ranging from "conventional" hardware to highly specialized silicon, optical, and molecular devices.

THE VECTOR MICROPROCESSOR SYSTEM

Vector processors, unlike conventional scalar processors, can specify multiple independent operations on linear operand arrays in one instruction. Thus, vector microprocessor architectures make ideal processing elements for multimedia and human machine interface. Applications in this area often contain algorithms that can be expressed in data-parallel form. A vector instruction-set architecture (ISA) allows a natural expression of the application's data parallelism. This simplifies implementations that adopt multiple parallel units and pipelined functional units. A well known vector microprocessor is the **TORRENT (TO)** system shown in Fig. 2. The TO processor is a complete single-chip Torrent architecture implementation fabricated by Hewlett-Packard's CMOS26G process. TO's main components are the MIPS-II-compatible RISC CPU with 1-Kbyte on-chip instruction cache, a vector unit coprocessor, an external memory interface, and an 8-bit wide serial host interface (TSIP) and control unit.

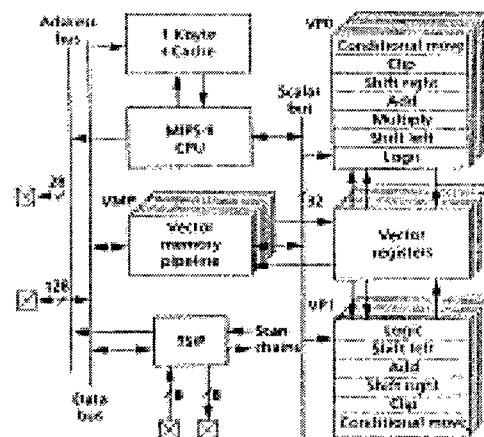


Fig. 2. Block diagram of TO micro architecture.

The latest development of vector processing system is the SPERT-II accelerator. It is a double-slot Sbus card for SUN-compatible workstations [12]. As shown in Fig. 3, the board contains a TO vector microprocessor, SRAM, a Xilinx FPGA device for interfacing with the host, and various system support devices.

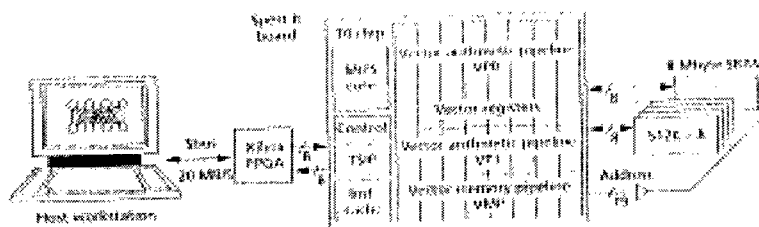


Fig. 3. Spert-II system organization.

SPERT-II software environment is shown in Fig. 4. The unmodified version of the gcc C/C++ compiler operates in parallel with the Torrent vector instructions which act as coprocessor instructions to the base MIPS-II instruction set.

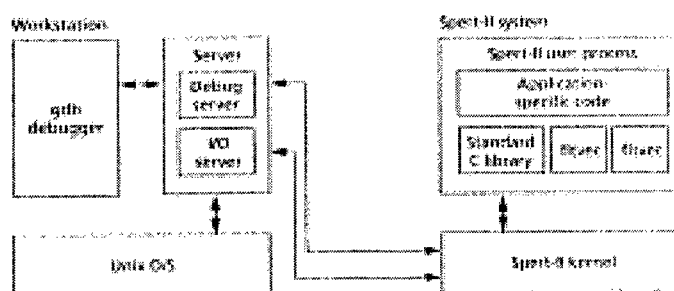


Fig. 4. The Spert-II software environment.

As described in [13], an (ANN) training task from a speaker-independent continuous speech recognition system was implemented on such SPERT-II system. The ANN is a simple three-layer, feed-forward perceptron with 100 to 400 input units, input layer fully connected to a hidden layer of 100 to 4000 units which, in turn, is fully connected to an output layer of 56 to 61 units. The back propagation algorithm was mapped with the hidden unit incorporating the standard sigmoid nonlinearity, while the output units compute a soft-max $f(x) = e^x / \sum S_i e^x$. The authors reported relative fast convergence of the algorithm: a randomly initialized net containing from 40,000 to 1 million weights converged with three passes over a training database containing several million patterns representations. A comparative performance is illustrated in Table 3.

Cellular Neural Networks

Cellular Neural Networks (CNN), proposed in 1988 by L. O. Chua and L. Yang [14, 15], are examples of recurrent networks defined by the following system of differential equations:

$$\frac{dx_{ij}(t)}{dt} = -x_{ij}(t) + \sum_{mn \in N_{ij}} a_{mn} y_{mn}(t) + \sum_{mn \in N_{ij}} b_{mn} u_{mn} + I \quad 10.$$

where

N_{ij} denotes the neighborhood of the ij -th cell for $1 \leq i \leq M$, $1 \leq j \leq N$ and $y = (|x + 1|) - (|x - 1|)/2$.

The state input and output of a cell are defined by x_{ij} , u_{ij} , and y_{ij} , respectively. The nearest neighborhood CNN is assumed. The output at an equilibrium point, if one exists, is denoted by y_{ij} . The parameters of a CNN are gathered into the so-called A-template, and B-template and the bias I . Since a CNN is a recurrent neural network, one can apply a suitable learning algorithm such as the recurrent back propagation (RBP) by replacing the standard sigmoidal function by a smoother similar one. However, there are still some unresolved issues concerning the minimization of the cost function associated with this type of network [16]. Nevertheless, CNNs are of great interest due to the fact that such networks are among the easiest to implement in hardware. In fact, CNNs are a particular class of artificial neural networks well suited for VLSI implementation.

Essentially, the CNN hardware structure consists of a bi-dimensional array of elementary analogue processors (cellular cells) locally interconnected only. Recently, the literature has reported several hardware implementations in terms of circuitry as well as the level of programmability. The Digitally Programmable CNN (DPCNN) is the latest of a neural chip family whose main feature is represented by a multi chip approach [17]. This feature, realized by using a current-mode approach, allows to implement any size CNN arrays by simply connecting many of these chips together. Moreover, a local digital memory implemented on each chip, enables the digital programmability of the template entries. Currently, DPCNNs have been applied to analog encryption and secure communications, analog built-in self-test, stochastic neural networks, annealing optimization and learning, autowaves for motion control, and biological inspired walking robots [18, 19].

CONCLUSION

Developments in ANNs have stimulated considerable enthusiasm and criticism. Some comparative studies are optimistic, some offer pessimism. The choice of the best technique should be driven by the nature of the application. Applications such as pattern recognition, clustering and categorization, function approximation, prediction/forecasting, optimization, retrieval by content, control, and plant process control, seem to be better served by combining the strengths of ANNs with other technologies. Such hybrid systems tend to achieve significantly better performance for these challenging problems. Vector microprocessors and VLSI technology of CNNs seem to be an optimistic technology for many new applications, such as those in *multimedia* and other *human-machine* interfaces. It is clear that future efforts will include communication and cooperative work between researchers working in ANNs and other disciplines so that repetitious work will be avoided, while good and effective methodologies for integration will be developed.

REFERENCES

1. J. Feldman, M.A. Fanty, N.H. Goddard, 1988. "Computing with Structured Neural Networks", Computer, 21(3), 91-103.
2. W.S. McCulloch, W. Pitts, 1943. "A Logical Calculus of Ideas Immanent in Nervous Activity," Bull. Math. Biophys., 5, 115-133.
3. D.O. Hebb, 1949. The Organization of Behavior, John Wiley & Sons, New York.
4. M. Minsky, S. Papert, 1969. Perceptrons: an Introduction to Computational Geometry, MIT Press, Cambridge, MA.
5. S. Haykin, Neural Networks: A Comprehensive Foundation, McMillan College Publishing Co., New York, 1994.
6. T. Kohonen, Self Organization and Associative Memory, Third Edition, Springer-Verlag, NY, 1989.
7. J. Hertz, A. Krogh, R.G. Palmer, 1991. Introduction to the Theory of Neural Computation, Addison-Wesley, MA.
8. G. A. Carpenter, S. Grossberg, 1991. Pattern Recognition by Self-Organization Neural Networks, MIT Press, MA.
9. J.J. Hopfield, 1982. "Neural Networks and Physical Systems with Emergent Collective Computational Abilities," in Proc. Nat'l Academy of Sciences, USA 79, 2,554 - 2,558.
10. Y. Shang, B.W. Wash, 1996. "Global Optimization for Neural Network Training", Comp., 28(3), 45-54.
11. T. Nordstrom, B. Svensson, 1992. "Using and Designing Massively Parallel Computers for Artificial Neural Networks," J. Parallel and Distributed Computing, 14(3), 260 - 285.
12. J. Wawrzyniec, K. Asanovic, N. Morgan, 1993. "The Design of a Neuromicroprocessor," IEEE Trans. Neural Networks, 4(3), 394-399.
13. J. Wawrzyniec et al., 1997. "Spart-II: A vector Microprocessor System," Computer, 29(3), 81-85.
14. L.O. Chua, L. Yang, 1988. "Cellular Neural Networks: Theory and Applications," IEEE Trans. Circuit and Systems, 32, 1257 - 1290.
15. B. Mirzai, D. Lim, G.S. Mschytz, 1996. "Robust CNN Templates: Theory and Simulations," IEEE. Workshop on Cellular Neural Networks and their Applications, Seville, 393-398.
16. B. Mirzai, Z. Cheng, G.S. Mschytz, 1998. "Learning Algorithms For Cellular Neural Networks," Proc. IEEE Int. Conf. On Circuits and Systems, (ISCAS-98), 3, Monterey, Ca. USA, 159-162.

17. M. Salerno, F. Sargeni, V. Bonaiuto, 1997. "A 720 Cells Interconnection-Oriented System for Cellular Neural Networks," Proc. IEEE Int. Conf. On Circuits & Systems (ISCAS-97), Hong Kong, 681-684.
18. _____, "An Improved Architecture for the Interconnections in a Multi-Chip CNN System," Proc. Int. Conf. On Circuits and Systems, (ISCAS-98), 3, Monterey, Ca. USA, 143-146.
19. P. Arena, L. Fortuna, M. Branciforte, P. Di Grazia, "Autowave for Motion Control: A CNN Approach," Proc. Int. Conf. On Circuits and Systems, (ISCAS-98), 3, Monterey, Ca., USA, 163-166.

Logical Rule Extraction from Data by Maximum Neural Networks

T. Saito and Y. Takefuji

Graduate School of Media and Governance,
Keio University, Fujisawa, Kanagawa, Japan

Abstract

In this paper, a new neural computing method to extract logical rules from the training data sets is proposed. Maximum neural networks are used to train the weight and the threshold of the multi-layered (feedforward) neural network (MLNN). The threshold and the weights of the MLNN are trained to be a logical function (AND/OR) with the multiple input. The maximum neural network constructs the logical function on the MLNN so that it is not necessary to extract rules from the trained MLNN. The proposed method was experimented for the classification problem, Monk's problem 1. Experimental results showed that the proposed method learned the correct rule in more than 40% success rate.

INTRODUCTION

The multi-layered (feedforward) neural network (MLNN) has been used to analyze data and has shown its effectiveness in several fields where Back-propagation is used as a learning method to train the weights of the MLNN. However, it is difficult to explain what the MLNN has learned by only seeing the weights of the network as information in the training data set is distributed into the weights and expressed across the whole network. So, rules must be extracted from the trained MLNN to understand what the MLNN has learned from the data. Several approaches to extracting rules from trained MLNNs have been suggested. Fu proposed the KT algorithm [1]. Towell et al. developed the M-of-N algorithm [2]. Ishikawa proposed a structural learning method [3]. Duch et al. proposed a method to extract logical rules for a Back-propagation MLNN. [4] However, extracting rules from a trained MLNN is very complex. It is shown that extracting DNF (Disjunctive Normal Form) rules from a trained MLNN is an NP-Hard problem [5].

This paper presents a new learning method for MLNN to extract logical rules from the training data using a maximum neural network. In the proposed method, one of three values $\{-1, 0, 1\}$ is assigned as a weight to the synaptic links of the MLNN. The neurons in the hidden and output layers are AND/OR functions with multiple input. The threshold value of each neuron is calculated according to the mode of the function (AND/OR) and the weights of the synaptic links with neurons in the lower layer. So the MLNN consists of logical functions. Rules can be obtained from the weights and threshold values without analysis. Therefore, extracting rules from the trained MLNN is unnecessary. The maximum neural network is used to train the network and select the mode of the function. The proposed method can be used for binary classification problems to extract DNF or CNF (Conjunctive Normal Form) rules. The number of output neurons is one. The data treated in the proposed method are binary but symbolic inputs are required to create logical rules

The search space of this problem is $3^N \times 2^M$ where N is the number of synaptic links and M is the number of neurons in the MLNN. When the problem size is large, it is difficult to find an optimal solution (the simplest rule) in practical time. The proposed method uses a heuristic search by a maximum neural network to find the (semi-)optimal solution in practical time.

MLNN

The MLNN used in the proposed method consists of three layers: input, hidden and output layers. The MLNN and the data sets have the following characteristics

- The data sets consist of the input data x_{pj} and the target output data t_p where p denotes the p th pattern and j denotes the j th attribute of the input pattern. x_{pj} is the binary data: the value 1 means "true" and 0 means "false" or "not". t_p has a binary value: 1 for positive examples and 0 for negative examples

- Each synaptic link has one of the three value of $\{-1, 0, 1\}$ as a weight value. The weight value 1 means “true”, -1 means “not” and 0 means “cut link”.
- Each neuron in the hidden and output layers consists of McCulloch-Pitts binary neuron which is formalized as follows,

$$o_j = \begin{cases} 1, & \text{if } \sum_{i=1}^C w_{j,i} o_i > \theta \\ 0, & \text{otherwise} \end{cases} \quad 1.$$

where o_j is the output of neuron in the current layer, o_i is the output of the neuron in the lower layer, $w_{j,i}$ is the weight of the synaptic link from i th neuron to j th neuron, C is the number of the synaptic links connected to j th neuron and θ is the threshold. The output 1 means “true” and 0 means “false” or “not”.

- The threshold value for neurons in the hidden and output layers can be calculated using weights of synaptic links from the lower layer according to the function mode:

$$\theta = \begin{cases} \sum_{w_{j,i} \in Z^+} w_{j,i} - 0.5, & \text{if the logical function mode is AND.} \\ \sum_{w_{j,i} \in Z^-} w_{j,i} + 0.5, & \text{if the logical function mode is OR} \end{cases} \quad 2.$$

where Z^+ is the set of weights which have the value “1” and Z^- is the set of weights which have the value “-1”. This definition is based on KBANN-net proposed in [2].

In this paper, one and only one output neuron is used to solve the binary classification problem. The number of hidden neurons can be changed according to problems. Neurons in the MLNN are numbered from 1 to I for the input neurons, from $I+1$ to $I+H$ for the hidden neurons and $I+H+O(=M)$ for the output neuron, where I , H and O are the numbers of neurons in the input, hidden and output layers respectively (Fig.1(a) shows an example).

The objective of this problem is to maximize the following function:

$$R = \left(1 - \frac{E_{pos}}{N_{pos}}\right) + \left(1 - \frac{E_{neg}}{N_{neg}}\right) \quad \text{where } E_{pos} = \sum_{p \in T_{pos}} |t^p - o_M^p|, \quad E_{neg} = \sum_{p \in T_{neg}} |t^p - o_M^p| \quad 3.$$

t^p and o_M^p are the target output and the output from the MLNN for the input pattern respectively. N_{pos} and N_{neg} are the numbers of positive and negative examples in the data respectively. T_{pos} and T_{neg} are the sets of the target outputs in the positive and negative examples respectively. R is 2 when the positive examples and the negative examples are learned correctly.

MAXIMUM NEURON

The maximum neuron consists of several neurons [6]. One and only one neuron in a maximum neuron can be activated. The definition of the maximum neuron is give by:

$$V_{m,n} = \begin{cases} 1, & U_{m,n} = \max (U_{x,n}) \quad x=1, \dots, X \\ 0, & \text{otherwise} \end{cases} \quad 4.$$

where $V_{m,n}$ and $U_{m,n}$ are the input and the output of m , n th neuron respectively and X is the number of neurons in one maximum neuron. In the proposed method, the difference between the maximum U and the minimum U is limited within D where D is a constant parameter.

NEURAL REPRESENTATION

In the proposed method, two maximum neural networks are used: one is for selecting the weight value and another is for determining the logical function mode. Fig.1 shows the MLNN and the neural representations corresponding to the MLNN. Fig.1(a) is the MLNN which has three input neurons, two hidden neurons and one output neuron. Fig.1(b) shows the neural representation of the maximum neural network for the weight value selection. This is represented by $3 \times S$ matrix where S is the number of synaptic links in the MLNN.

Each column represents one maximum neuron so that one and only one neuron can be activated per a column. In this paper, first, second and third rows represent the weight value 1, -1, and 0 respectively. The black square means that the m,n th neuron in the maximum neural network is activated ($V_{m,n}=1$), in other words, this means the m th weight value is selected for the weight of the n th synaptic link. For example, in Fig.1(b) the state of 2,4th square is activated. This means that the weight value of the 4th synaptic link is -1. The white square means m,n th neuron is not activated ($V_{m,n}=0$). In the same way, Fig.1(c) shows the neural representation for selecting the logical function mode of the neuron. This representation is $2 \times L$ matrix where L is the number of the neurons in the hidden and output layers. In this paper, the first and second rows represent AND function and OR function respectively.

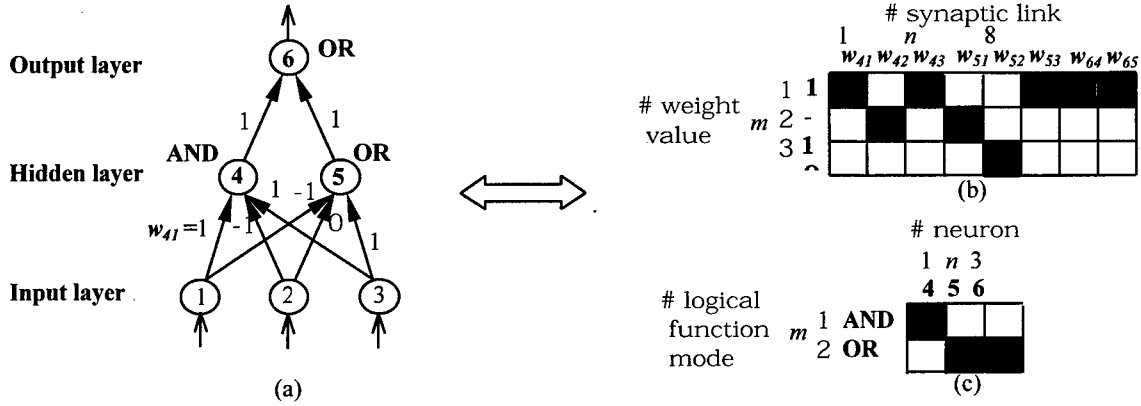


Fig.1: (a)MLNN and neural representations (b) for the weight selection and (c) for the logical function mode selection.

MOTION EQUATION

The input U of the maximum neuron is updated by the motion equation. The motion equation is given by:

$$\frac{dU_{m,n}}{dt} = -\alpha \left\{ 3 - \left(\left(1 - \frac{E_{pos}}{N_{pos}} \right) + \left(1 - \frac{E_{neg}}{N_{neg}} \right) + \left(1 - \frac{E_{pos}}{N_{pos}} \right) \times \left(1 - \frac{E_{neg}}{N_{neg}} \right) \right) \right\} \quad 5.$$

α is a constant parameter. When the positive examples and the negative examples are learned correctly, Eq.5 becomes zero. This motion equation is used for two maximum neural networks. In this motion equation, there is an objective function. The method where the objective function is used as a motion equation has been proposed by Tsuchiya et al. [7]. In the proposed method, in the case that n represents the weight from the i th input neuron to the j th hidden neuron (w_{ji}) and the weight value from the j th neuron to the output neuron is zero ($w_{Mj}=0$), either the weight value 1 or -1 which has larger U in the k th maximum neuron corresponding to $w_{M,k}$ is set to $w_{M,j}$ and then o^p is calculated. The same operation is applied to the calculation of the motion equation for the n th hidden neuron

The following term is added to Eq.5 to escape from the state of the tautology:

$$-\beta \sum_{\substack{I+1 \leq k \leq I+H, \\ j \neq k}} f_1(w_{ji}, w_{ki}) \quad 6.$$

where β is the constant parameter and w_{ji} is the weight of the n th synaptic link. This term is applied to the calculation of the motion equation for the synaptic links from the input neuron to the hidden neuron ($1 \leq i \leq I$). The function f_1 returns 1 if the state of the tautology is detected, otherwise 0. Let r_i denote the logical function mode of the i th neuron. The state of the tautology related to the n th synaptic link, w_{ji} , can be detected by using the following rules:

Under the condition that $w_{ji} = w_{ki} \neq 0$ and $w_{ji} \neq w_{ki}$,

1. $r_i = r_k = r_M = \text{OR}$, $w_{M,j} = w_{M,k} = 1$
2. $r_i = r_k = \text{OR}$, $r_M = \text{AND}$, $w_{M,j} = w_{M,k} = -1$
3. $r_i = r_k = \text{AND}$, $r_M = \text{OR}$, $w_{M,j} = w_{M,k} = -1$
4. $r_i = r_k = r_M = \text{AND}$, $w_{M,j} = w_{M,k} = 1$

or under the condition that $w_{j,i} = w_{k,i} \neq 0$ and $w_{j,i} = w_{k,i}$,

5. $r_j = r_M = \text{OR}$, $r_k = \text{AND}$, $w_{M,j} = 1$, $w_{M,k} = 1$
6. $r_j = \text{OR}$, $r_M = r_k = \text{AND}$, $w_{M,j} = -1$, $w_{M,k} = 1$
7. $r_j = \text{AND}$, $r_M = r_k = \text{OR}$, $w_{M,j} = -1$, $w_{M,k} = 1$
8. $r_i = r_M = \text{AND}$, $r_k = \text{OR}$, $w_{M,j} = 1$, $w_{M,k} = -1$

or under the condition that $w_{j,i} = w_{k,i} \neq 0$ and $w_{j,i} \neq w_{k,i}$ and $\sum_{1 \leq i \leq I} f_2(w_{j,i}) = 1$ and $\sum_{1 \leq i \leq I} f_2(w_{k,i}) = 1$,

9. $r_i = r_k = \text{AND}$, $r_M = \text{OR}$, $w_{M,j} = w_{M,k} = 1$
10. $r_i = r_k = r_M = \text{OR}$, $w_{M,j} = w_{M,k} = -1$
11. $r_i = r_k = \text{OR}$, $r_M = \text{AND}$, $w_{M,j} = w_{M,k} = 1$
12. $r_i = r_k = r_M = \text{AND}$, $w_{M,j} = w_{M,k} = -1$

or under the condition that $w_{j,i} = w_{k,i} \neq 0$ and $w_{j,i} \neq w_{k,i}$ and $\sum_{1 \leq i \leq I} f_2(w_{j,i}) = 1$ and $\sum_{1 \leq i \leq I} f_2(w_{k,i}) = 1$

13. $r_j = \text{AND}$, $r_M = r_k = \text{OR}$, $w_{M,j} = 1$, $w_{M,k} = -1$
14. $r_j = r_M = \text{OR}$, $r_k = \text{AND}$, $w_{M,j} = -1$, $w_{M,k} = 1$
15. $r_j = \text{OR}$, $r_M = r_k = \text{AND}$, $w_{M,j} = 1$, $w_{M,k} = -1$
16. $r_i = r_M = \text{AND}$, $r_k = \text{OR}$, $w_{M,j} = -1$, $w_{M,k} = 1$

where the function f_2 returns 1 if the weight value is not 0, otherwise 0.

To escape from the local minimum, the following motion equation is used instead of Eq.5:

$$\frac{dU_{m,n}}{dt} = -\alpha \left\{ 3 - \left(\left(1 - \frac{E_{pos}}{N_{pos}} \right) + \left(1 - \frac{E_{neg}}{N_{neg}} \right) + \left(1 - \frac{E_{pos}}{N_{pos}} \right) \times \left(1 - \frac{E_{neg}}{N_{neg}} \right) \right) \right\} V_{m,n} \quad 7.$$

Eq.7 is applied to the neuron in the maximum neural network whose output is 1. This term is used if t modulo $a > b$ where t is the iteration step and a, b are constants.

ALGORITHM

The following procedure describes the proposed algorithm. Note that t_limit is the maximum number of iteration steps for the system termination condition, $target_r$ is the target value of R and, $UW_{m,n}$, $VW_{m,n}$ and $UL_{m,n}$, $VL_{m,n}$ denote the input U and the output V of two maximum neural networks for the weight selection and the logical function selection respectively.

1. Set $t=0$ and set $\alpha, \beta, \gamma, t_limit, D, a, b$.
2. The initial values of $UW_{m,n}(t)$ for $m, \dots, 3, n, \dots, N$ and $UL_{m,n}(t)$ for $m, \dots, 2, n, \dots, H+O$ are randomized.
3. Evaluate $VW_{m,n}$ for $m, \dots, 3, n, \dots, N$ and $VL_{m,n}$ for $m, \dots, 2, n, \dots, H+O$ using Eq.3.
4. Set the weight values to the MLNN and calculate the thresholds. If $R \geq target_r$, then terminate this procedure and go to step 8.
5. For $n=1, \dots, N$,
 - (a) Set the m th weight value to the n th synaptic link and compute $\Delta UW_{m,n}(t)$ using Eq.5 and Eq.6 for $m=1, \dots, 3$. If R in Eq.5 becomes larger than $target_r$, then terminate this process.
 - (b) If the n th synaptic link represents $w_{j,i}$ ($1 \leq i \leq I, I+1 \leq j \leq I+H$) and the value of Eq.5 of $\Delta UW_{1,n}(t)$ or $\Delta UW_{2,n}(t)$ are equal to the value of Eq.5 of $\Delta UW_{3,n}(t)$, then the following equation is applied to $\Delta UW_{m,n}(t)$ ($m=1$ or 2):

$$-\gamma \left(\sum_{1 \leq i \leq I, I+1 \leq j \leq I+H} f_2(w_{j,i}) + \sum_{I+1 \leq j \leq I+H, I+H+1 \leq k \leq M} f_2(w_{k,j}) \right) \quad 8.$$

where f_2 returns 1 if the weight value is not 0, otherwise 0.

- (c) Update $UW_{m,n}(t+1)$ for $m=1, \dots, 3$: $UW_{m,n}(t+1) = UW_{m,n}(t) + \Delta UW_{m,n}(t)$.
- (d) Evaluate $VW_{m,n}(t+1)$ for $m=1, \dots, 3$, using Eq.4.
6. The same procedure except 5(b) are applied to $UL_{m,n}$, $VL_{m,n}$.
7. If $t=t_limit$ then terminate this procedure else increment t by 1 and go to step 5.

8. Prune synaptic links whose weight values are 1 or -1 but does not affect the value of the objective function by setting the weight value to 0.

EXPERIMENTAL RESULTS

The proposed method was experimented with the binary classification benchmark problem, Monk's problem 1 (M_1). M_1 has six different attributes to describe an artificial robot domain. There are four rules to be a monk in M_1 . Table 1 describes the six attributes and rules. The binary strings in Table 1 are binary expressions of the attribute values. The input pattern consists of 15 binary values.

Table 1. : Attributes and Rules of Monk's Problem 1.

attributes	Values	Rules
head_shape	round(100), square(010), octagon(001)	head_shape=round and body_shape_round
body_shape	round(100), square(010), octagon(001)	head_shape=square and body_shape_square
is_smiling	yes/no (1/0)	head_shape=octagon and body_shape_octagon
holding	sword(100), balloon(010), flag(001)	jacket_color=red
jacket_color	red(1000), yellow(0100), green(0010), blue(0001)	
has_tie	yes/no (1/0)	

The proposed method was experimented a thousand times with different random numbers on PentiumII (450MHz) computer. In experiments, 15 input neurons, 5 hidden neurons and 1 output neuron were used. All input neurons were connected to 5 hidden neurons, which themselves were connected to the output neuron. $\alpha, \beta, \gamma, t_limit, D, a, b$ are set to 1.5, 0.1, 0.01, 2000, 1.5, 15, 13 respectively. The way of mapping the synaptic link in the MLNN to the maximum neuron is the same with one shown in Fig.1. The number of positive examples and negative examples are 64 respectively.

Fig.2 shows two MLNNs obtained in experiments. Table 3 shows the rules obtained from Fig.2(a) and Fig.2(b). In Fig.2(a), the logical function mode of the output neuron is OR and one of hidden neurons is AND. This means that the rule is expressed in DNF. In Fig.2(b), the logical function mode of the output neuron is AND and one of hidden neurons is OR. This means that the rule is expressed in CNF. Table 3 shows the computation result of 1000 times experiments.

Table 2: Rules obtained in experiments (corresponding to Fig.2(a) and Fig.2(b)).

Rule in Fig.2(a)	Rule in Fig.2(b)
(Jacket_color=red) or (head_shape=round and body_shape=round) or (head_shape \neq round and head_shape \neq octagon and body_shape=square) or (head_shape=octagon and body_shape=octagon)	(head_shape \neq round or body_shape=round or jacket_color=red) and (head_shape=round or head_shape=octagon or body_shape=square or jacket_color=red) and (head_shape \neq octagon or body_shape=octagon or jacket_color=red)

Table 3.: The result of the 1000 experiments: the average iteration steps, the average computation time (CPU time) and the convergence rate to the correct rule within 2000 iteration steps.

Average iteration steps	Average computation (CPU)time (sec.)	Convergence rate
806.1940	64.006	0.415

DISCUSSION AND CONCLUSION

In this paper, the logical rule extraction method with maximum neural networks is proposed. The proposed method could find several rules expressed in DNF and CNF for the monk's problem 1. Rules were extracted directly from the trained MLNN

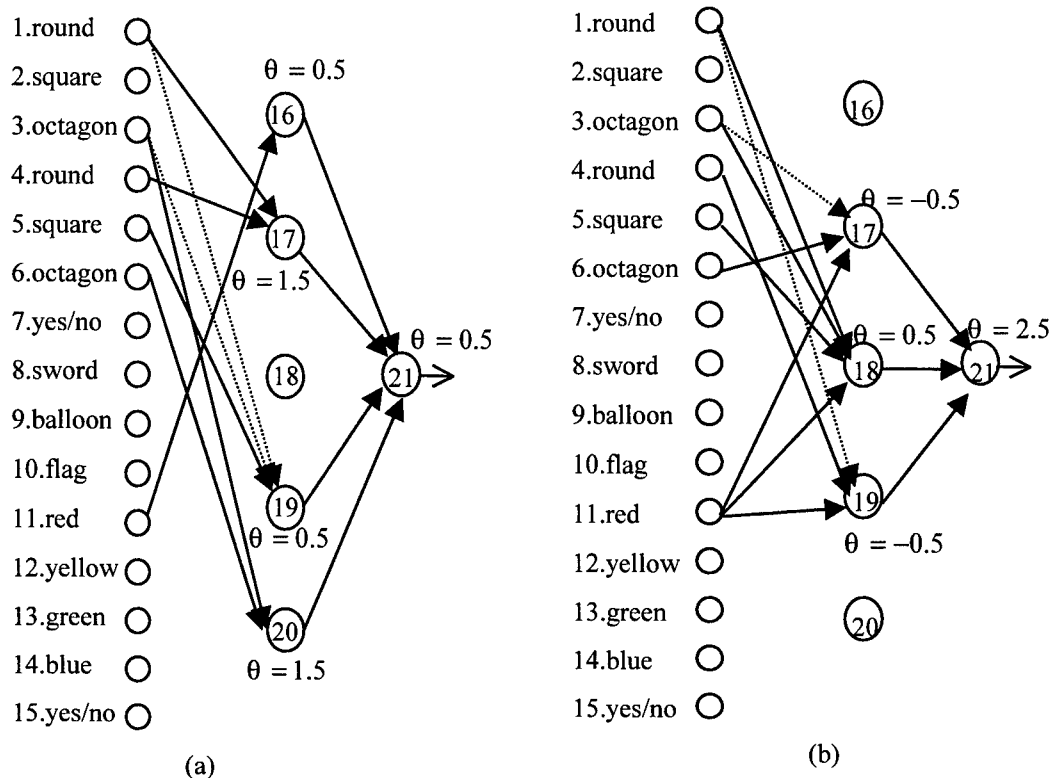


Fig.2: MLNNs obtained in experiments: an arrow denote the synaptic link whose weight value is 1 and the dashed arrow denotes the synaptic link whose weight value is -1. The synaptic link whose weight value is 0 is not displayed.

without analyzing weight values in the MLNN. However, the training data of the monks problem 1 does not contain the noise and the rule is comparatively simple. More experiments with the data including the noise and the complicated rules are needed to evaluate the performance of the proposed method. The computation time of the proposed method depends on the data size and the complexity of the rule. When the problem size is large, it takes much time to calculate the motion equation. The complicated rule requires many hidden neurons. At the same time, many synaptic links are required. Especially, the computation time for the motion equation should be reduced. These are future works to be solved for the improvement.

REFERENCES

1. L. M. Fu., 1994. "Rule Generation form Neural Networks", IEEE Transactions on Systems, Man and Cybernetics, 24(8), 1114-1124.
2. G. G. Towell, 1993. "Extracting Refined Rules from Knowledge-Based Neural Networks", Machine Learning, 13, 71-101.
3. M. Ishikawa, 1990. "A Structural Connectionist Learning Algorithm with Forgetting", J. of Japanese Society for Artificial Intelligence, 5(5), 595-603.
4. W. Duch, R. Adamczak, K. Grabczewski, 1996. "Extraction of logical rules from training data using backpropagation networks", 1st Online Workshop on Soft Computing
<<http://www.bioele.nuee.nagoya-u.ac.jp/wsc1/>>.
5. A.B.Tickle, M. Golea, R. Hayward, J. Diederich, 1997. "The Truth Is in There: Current Issues in Extracting Rules from Trained Feedforward Artificial Neural Networks", IEEE ICNN, 2530-2534.
6. Y. Takefuji, 1992. "Neural Network Parallel Computing," Kluwer Academic Publishers.
7. K. Tsuchiya, S. Bhariktar, Y. Takefuji, 1996. "A neural network approach to facility layout problem", European Journal of Operational Research, 89, 556-563.

Iterative RBF Neural Networks as Metamodels of Stochastic Simulations

George Meghabghab*, George Nasr**

* Valdosta State University, Department of Math & Computer Science,
Valdosta, GA 31698

Email: gmeghab@valdosta.edu

** Lebanese American University, Department of Electrical and Computer Engineering,
PO Box 36, Byblos, Lebanon, 13-5053.

Email: genasr@lau.edu.lb

ABSTRACT

Research into emerging technological approaches to make computer simulations more effective and efficient is an essential ingredient to developing successful manufacturing models. This study is a premiere study in using neural networks in metamodeling stochastic simulation in manufacturing domain. A new iterative RBF neural network was developed rather than the baseline ANN models which were used in stochastic simulation metamodeling in domains such as combat simulations in the military, service industries, and transportation companies. Given the fact that typical stochastic simulation metamodeling approaches involves the use of regression models in response surface methods, RBF become a natural target for such an attempt because they use a family of surfaces each of which naturally divides an input space into 2 regions and the n patterns will be assigned either class $X+$ or $X-$. This dichotomy of the points is said to be separable with respect to the family of surfaces if there exists a surface in the family that separates the points in the class $X+$ from those in the class $X-$. In fact, for the evaluation of the quality of a ball steel, RBF metamodel trained on 1521 training examples from a set of 13000 different simulation runs and was able to outperform direct simulation on 120 additional test examples which were not included in the training set.

INTRODUCTION

Computer simulations are widely used in a variety of applications including the military, service industries, manufacturing companies, nuclear power plants, and transportation organizations. For example in nuclear power plants, often computer simulations are used to train personnel on failures and normal operation, study operation plans, support testing of the heart of the nuclear power plant, as well as to evaluate and compare future design plan changes. Other techniques that are used to examine systems in general do not have the advantages that computer simulations bring mainly that computer simulations provide cheaper and realistic results than other approaches do. In some cases, computer simulation is the only means to examine a system like in nuclear power plants since it is too dangerous to bring a system under such failure conditions to study it closely or costly and infrequent like in combat situations to experiment with the system. Also computer simulation permit studying systems over large periods of time, learn from real world past experiences, and have control over experimental conditions.

Actual manufacturing simulation models are expensive to develop and use, in terms of personnel, time and resources. Large memory requirements, slow response time can prevent companies from considering it as a useful tool. The need to develop manufacturing simulations models that can be used in training that are as realistic as possible is the issue, and speed is not important, while in testing speed and reproducibility become important, incite us to make the different internal simulation modules as efficient and accurate as possible.

Computer simulations have provided companies with the description of the input settings that are needed to produce the optimal best output value for a given set of inputs in a specific domain of study. Response

surface methodologies using regression models approximations of the computer simulation were the means to achieve computer simulation optimization. As Myers, Khuri, and Carter (1989) stated it in *Technometrics* [1] that: "There is a need to develop non-parametric techniques in-response surface methodologies. The use of model-free techniques would avoid the assumption of model accuracy or low-level polynomial approximations and in, particular, the imposed symmetry, associated with a second degree polynomial". One possible non-parametric approach is to use artificial neural networks.

ANNs IN MANUFACTURING AND STEEL PRODUCTION

An Artificial Neural Network (ANN) learns to imitate the behavior of a complex function by being presented with examples of the inputs and outputs of the function. This operation is called training. A successfully trained ANN can predict the correct output of a function when presented with the input (or inputs for a multivariate function). ANNs have found wide applications in a variety of fields including mining and manufacturing.

Because there are many steel production processes that are complex and uncertain when modeling for control is concerned, there are examples of ANN being utilized in the steel industry. Neural network controllers have been developed for electric arc furnaces [2], for a continuous casting process [3], and the modeling of a quality steel production with an adaptive logic network [4]. In a totally different and new perspective, ANN were used to determine the surface glossiness of steel sheets as an evaluation method [5], while Lusiak and Pietrzyk [6] used ANN as a history dependent constitutive model for hot forming of steels. In the next section, we will examine the role of ANN in approximating computer simulations in a manufacturing domain mainly the evaluation of a quality steel production.

ARTIFICIAL NEURAL NETWORK METAMODEL APPROACH

A metamodel is a model of a model ([7]). Typical simulation metamodeling approaches involve the use of regression models in response surface methods. A recent overview of published research on simulation metamodels can be found in [8]. A few attempts have been made to employ neural networks as the metamodeling technique. Using an ANN to model a stochastic simulation was done [9], [10], [11], and [12]. Each of these researchers was successful in using ANN as metamodels of stochastic computer simulations. The common feature of these models was an ANN baseline which involves using a backpropagation trained, multi-layer ANN to learn the relationship between the simulation inputs and outputs. The baseline ANN metamodel approach was developed on a (s,S) inventory computer simulation and also was applied to a larger application in the domain under consideration.

RBF has been developed now for a number of years. There is a resurgence in using RBF as a viable architecture to implement neural network solution to many problems. RBF neural networks are deterministic global non-linear minimization methods. These methods detect subregions not containing the global minimum and exclude them from further consideration. In general, this approach is useful for problems requiring solution with guaranteed accuracy. These are computationally very expensive. The mathematical basis for RBF networks is provided by Cover's Theorem [13] which states that a nonlinearly-separable pattern classification problem in high-dimension space is more likely to be linearly-separable than in low-dimensional space. This is the reason for choosing a high dimension for the hidden layer in the network. RBF uses a curve-fitting scheme to learn, i.e., learning is equivalent to finding a surface in a multi-dimensional space that represents a best fit for the training data ([14], [15]).

The approach considered here is a generalized RBF neural network where the number of nodes at the hidden layer is M , where M is smaller than the number of training patterns N . At the output layer, the linear weights associated and the position of the centers of the radial basis functions and the norm weighting matrix associated with the hidden layer are all unknown parameters that have to be learned. A supervised learning process using a gradient descent ([16]) procedure is implemented to adapt the position of the centers and their spreads (or widths) and the weights of the output layer. To initialize such a gradient descent GD procedure or CCSW we begin the search from a structured initial condition that limits the region of parameter space to be searched to

an already known useful area through using a standard pattern-classification method as an RBF network ([17]). The likelihood of converging to an undesirable local minimum in weight space is already reduced. Also a supervised learning process using interior point method IPM developed in ([18], [19]) is implemented to adapt the position of the centers and their spreads(or widths) and the weights of the output layer but which reduces the amount of computation compared to the one developed of GD in ([16]). A standard Gaussian classifier is used which assumes that each pattern in each class is drawn from a full Gaussian distribution.

An iterative RBF metamodel approach to approximating discrete event computer simulations was developed in order to develop accurate metamodels of computer simulations. An RBF neural network was used to learn the relationship between the simulation inputs and outputs. The iterative RBF metamodel approach starts with small training and testing sets, in terms of replications, and build RBF metamodels to use in performing factor screening to eliminate those input factors that do not appear to make much a difference on the simulation output. After eliminating the irrelevant factors, the baseline approach could be used on the remaining factors with substantial savings in total computer simulation runs. The iterative RBF neural network was developed and was applied on an offline evaluation of the quality of a ball steel production line providing grinding media for the mining industry.

RESULTS OF RBF TRAINING ALGORITHM IN MANUFACTURING DOMAIN

The purpose of this research is to predict the proportion of rejected bars after they have been tested for voids larger than the maximum acceptable size. During this process, 9 significant variables were chosen as inputs based on a priori knowledge. The output variable is the defective percentage on a given example. A sample or an example of a given cast is said to be rejected if its defective percentage exceeds 40%. This last cutoff number might sound too high but based on all the parameters considered during the testing by different engineers it represents an accurate information of what a bar that is accepted and a bar that is rejected. Our data were limited to 70 different examples because all inputs and output were complete for these different examples. The 9 different inputs are proprietary data for the steel company and cannot be made public. This poses problems for those that try to simulate our data on different regression models since the range of these different inputs cannot be released.

First order multiple linear regression models were applied to the different examples. The F test for regression was significant for $\alpha=0.01$. The regression model of the first order linear equation is of the form:

$$\text{Percentage Defective} = A + BX_1 + CX_2 + DX_3 + EX_4 + FX_5 + GX_6 + HX_7 + IX_8 + JX_9 \quad 1.$$

where A, B, C, D, E, F, G, H, I, J where constants that were identified during simulation, and $X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8, X_9$ represent the 9 input variables. This process is a complex one since different input variables have different range values and different sampling values between these different ranges had to be tested for the best fit among all these constants. If each input variable had just a maximum of 5 different sampling values, there are 1048576 (4^9) training sets possible for the different sets of values of the input parameters considered. Only few of these different were selected that range from 512 points minimum to 5^9 points maximum. Out of all of these possible training sets, 13000 were investigated but only one case is reported here.

Figure 1 shows the result of the actual output for these different examples. Out of the 70 examples, 60 were used on training and the remaining examples on testing. Figure 1 also shows the result of the first order linear regression of equation (1) applied to the different examples. RBF were trained on the 13000 different simulation cases. The model was made out of 9 different input neurons, one output neuron for each range of percentage defective, e.g. each range spans 10 % on the Y axis of Figure 1 and as such 7 neurons are needed for the seven ranges from 0% to 70%. The number of hidden layers depend on the number of exemplars from the 13000 cases for each range. There are more exemplars that cover a given range than other ranges. That makes the neurons learn more about a given range than other ranges. There were 1521 nodes at the hidden layer for that particular run. The success rate was 95.7% on all the examples for RBF without CCSW. RBF with CCSW using GD was 96.7% on all the examples. RBF with CCSW using IPM was up to 97.6% which is impressive compared to other studies done with RBF on other benchmark problems ([16], [18], [19]). The error

was of the order of 0.005. The matrices manipulations were quite heavy at the hidden layer level since a (1521*1521) matrix had to be inverted. RBF is heavy computationally. RBF without CCSW was of the order of 60 minutes. RBF with CCSW using GD was of the order of 120 minutes. RBF with CCSW using IPM was of the order of 90 minutes. The time is comparable with the first order linear regression analysis. From figure 1 we could see that RBF outperformed the first linear order regression analysis on all examples. Thus, RBF did very well on the training patterns. Since the set of testing is only limited to ten examples, the vector of weights were used but with less centers since the number of patterns in each range has shrunk considerable. On the testing patterns the success rate was of the order of 90% without CCSW, 91.2% with CCSW using GD, and 91.8% with CCSW using IPM. The results for testing were not shown but they follow the results on training.

The next step was to reconsider the high number of inputs that is needed in the process. Other studies([12]) suggested testing whether the actual metamodel supports the reduction of the number of inputs without affecting the performance of the output results specially when the number of inputs is high. This idea was tested effectively in our metamodel. Given the fact that we know the output of RBF with the complete set of 9 inputs, we could try eliminating in each run a given input and calculate the output of RBF and compared with the known output when the set of inputs is 9. If the error between both values was still smaller than a given threshold that particular input could be ignored. Thus 9 different runs were investigated each time eliminating a given input from the calculation. The results show that RBF was successful in eliminating 2 inputs out of 9 without affecting the performance of the RBF networks and were still better than first order linear regression analysis with 9 inputs. When we tried to combine eliminating more than 2 inputs at a time the network degraded enormously and the success rate dropped to less than 50% on the training patterns.

In addition, it was shown that networks trained on individual replication output data had better generalization performance than networks trained on only the averages of the simulation output. In other words it is best when approximating stochastic computer simulations to use "noisy" individual replications rather than the "quiet" average values. The iterative RBF metamodel performed well at approximating computer simulations. The iterative RBF metamodel approach can be used by other researchers for comparison purposes when developing their own RBF metamodel approach. This contribution is useful in the area of artificial neural networks because there are many different existing and emerging ANN procedures to perform approximation and estimation tasks.

BREAKTHROUGH ASPECT OF THE WORK AND CONCLUSION

A major area of research is in the experimental design of neural networks as metamodels of computer simulations. This research filled in a critical need in the designs that take into consideration both the development(training) and the evaluation (testing and validation) of metamodels. This research shows how a metamodel is to be constructed using a training set for adjusting the weights, the centers and the spreads, and one test set for determining when to stop training and a second test set for evaluation of the generalization ability of the metamodel which was never done before. This is a premiere study which uses an iterative approach rather than baseline ANN metamodeling ([12]) which is a major improvement in that it reduces simulation runs in almost 40%. This research investigated the use of extreme values observed from the simulation. The studies cited above ([9], [10], [11], [12]) ignored such values from the training set. Integrating values of the average output for a particular combination of the input parameters that are much larger than all other average output values was possible because of the nature of RBF neural networks. This research further enhances RBF neural networks as de facto neural networks when off line analysis is needed where speed is not the goal but accuracy is the final issue. It was shown on the 4 benchmark problems that RBF outperformed the best and the fastest technique in two problems out of 4 in training and testing and in testing only on another one. Also the number of original inputs in the evaluation of steel ball quality which was quite high, e.g. 9, was reduced during the study to seven which is major finding compared to other studies which suggested the reduction but failed to achieve it on their own data. The manufacturing domain represented a challenge to RBF and RBF faired better than other simulation analysis techniques and other backpropagation neural network techniques. Second order linear regression analysis techniques could have been used but the researchers believed that this would escalated the time to calculate the results and would not have improved the results of the first order linear regression analysis.

The authors of this study encourage other simulation analysts to use this study and our RBF as a model to build on in their actual domain and simulation. Research along these lines is essential to ensuring that this tool is properly integrated with other emerging technologies to provide successful future generations of manufacturing simulations which will save millions of dollars in quality production.

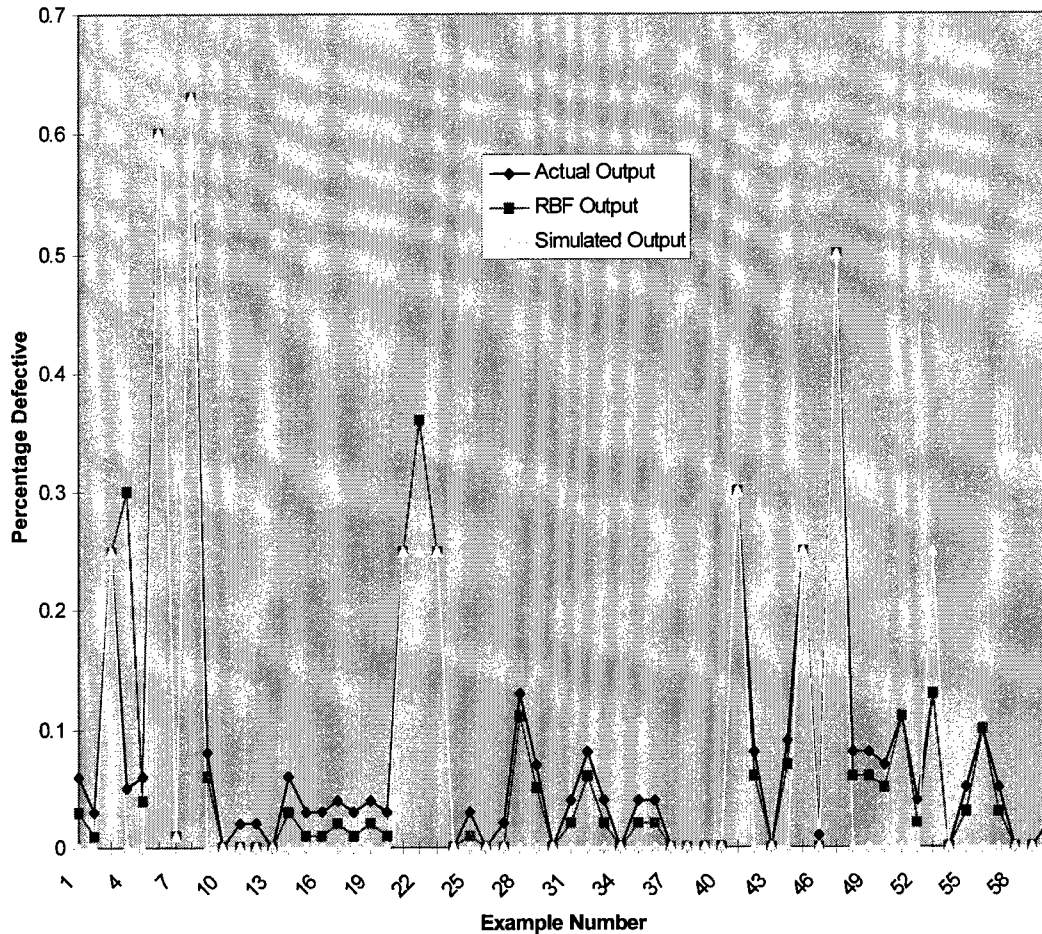


Fig. 1. Comparison of simulation results and ANN results with actual data.

REFERENCES

1. Myers, R., Khuri, A., and Carter W., 1989. Response Surface Methodology:1966-1988. *Technometrics*, 31(2), 137-157.
2. Staib, W.E, Bliss, N.G., and Staib, R.B., 1991. Neural Network Conversion of the Electric Arc Furnace Into the Intelligent Arc Furnace. *Iron and Steel Conference*, Washington D.C, April 15.
3. Kominami, H., Naitoh, S., Kamada, N., Hqamaguchi, C., Tanaka, T., and Endoh, H., 1991. Neural Network System for Breakout Prediction in Continuous Casting Process. Tech Rept. 49, Nippon Steel, 6-3, Otemachi 2-chrome, Chiyoda-ku, Tokyo 100-71, Japan, April.
4. Stelljes, T.A., and Erickson, K.T., 1995. Modeling the Quality of Steel Production With an Adaptive Logic Network. In *Intelligent Engineering Systems Through Artificial Neural Networks*, Vol.5., Fuzzy Logic and Evolutionary Programming, by Dagli et al. (Eds), ASME Press.
5. Tateno, J., Asano, K., Moriya S., and Shiokawa, T., 1997. Neural Network Based Evaluation Method for Surface Glossiness of Steel Sheets. In the *First International Conference on Intelligent Processing and Manufacturing of Materials*.

6. Kusiak, J. and Pietrzyk, M., 1997. Artificial Neural Networks as a History Dependent Constitutive Model for Forming of Steels. In the 1st International Conf. on Intelligent Processing and Manufacturing of Materials.
7. Blanning, R., 1975. The Construction and Implementation of Metamodels. *Simulation*, 24, 177-184.
8. Yu, B., and Popplewell, K. Metamodels in Manufacturing: A Review. *International Journal of Production Research*, 32, 787-796.
9. Pierreval, H., and Huntsinger, R., 1992. An Investigation on Neural Network Capabilities as Simulation Metamodels. *Proceedings of the 1992 Summer Computation Simulation Conference*, 413-417, Reno, Nevada, July 27-30.
10. Hurron, R.D., 1992. Using a Neural Network to Enhance the Decision Making Quality of a Visual Interactive Simulation Model. *Journal of The Operations Research Society*, 43(4), 333-341.
11. Badiru, A.B., and Sieger, D.B., 1993. Neural Network as a simulation metamodel in economic analysis of risky projects. Tech. Rept., Dept of Industrial Engineering, University of Oklahoma.
12. Kilmer, R.A., 1995. Applying Artificial Neural Networks to Combat Simulations. *Mathematical And Computer Modelling*, 23(1-2), 91-99.
13. T.M. Cover, 1965. Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. *IEEE Transactions on Electronic Computers*, EC-14, 326-334.
14. Powell, M.I.D., 1985. Radial Basis functions for multivariate interpolation: A review. *IMA Conference on Algorithms for the approximation of functions and data*, RMCS, Shrivenham.
15. Broomhead, D.S. and Lowe, D., 1988. Multivariable functional interpolation and adaptive networks. *Complex Systems*, Vol. 2, 321-355.
16. Meghabghab, G. , Nasr, G., and Boyd, D., 1997. Radial Basis Functions Neural Networks VS NOVEL on 4 benchmarks problems. *Proceedings of the 10th International FLAIRS Conference*, Daytona Beach, FL, 10-14 May 1997, 242-246.
17. Lowe, D, 1991. What have neural networks to offer statistical pattern processing. *Proceedings of the SPIE Conference on Adaptive Signal Processing*, 460-471, San Diego, CA.
18. Meghabghab, G. and Nasr, G. 1997. A new Radial Basis Function Neural Network VS NOVEL on 4 benchmarks problems. *Intelligent Engineering Through Artificial Neural Networks*, Edited by C. Dagli, Vol 7, ASME Press, New York, NY, 177-182.
19. Meghabghab, G. and Nasr, G. 1998. An Interior Point Radial Basis Function Neural Network on 4 benchmarks problems. *Proceedings of 4th World Congress on Expert Systems*, Mexico City, Mexico, March 16-20, 844-855.

A Systematic and Reliable Approach to Pattern Classification

R.Doraiswami, M.Stevenson, S. Rajan

The University of New Brunswick, Fredericton, NB. Canada E3B 5A3

ABSTRACT

A systematic and reliable approach to classify patterns is proposed when no *a priori* information except a set of pre-classified data is provided. A classifier is selected from a number of state of the art pattern classification schemes which are diverse in approach as well as the assumptions employed in their design. The selected schemes include the K-Nearest Neighbour Classifier (KNNC), the Minimum Mahalanobis Distance Classifier (MMDC), and the Artificial Neural Network Classifier (ANNC). In order to ensure that the selected classification scheme is properly designed and correctly implemented, the given pre-classified data is analysed, and the relative performance of the classifiers are cross validated as well as compared with a benchmark performance measure. The given data set is subjected to data validation, data visualization and feature quality analysis with a view to detect bad data, to obtain a qualitative picture of the class separability, and to derive a benchmark performance measure called the Bhattacharyya distance measure. In the design phase, the classifiers are executed in the order of increasing accuracy and increasing complexity so that a classifier at one level in the hierarchy sets the performance goal (e.g. classification accuracy) for the task at the next level. Further, to ensure a peak performance, the classifier accuracy is compared with the Bhattacharyya distance measure. The proposed scheme is evaluated on both simulated as well as actual data obtained from the images of the biological cells.

INTRODUCTION

Classification of patterns is a challenging task, and finds wide applications in many fields including character recognition, fingerprint classification, medical diagnosis, automatic target recognition, industrial inspection, machine vision, visual servoing, fault diagnosis, and speech recognition. Pattern classification includes identification of a set of features, computation of the features and classification using the features. See Fig.1.

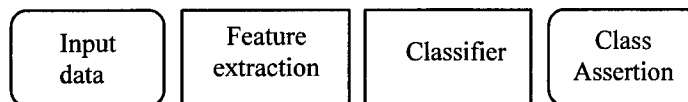


Fig. 1. Components of Pattern Classification

The success of the classifier depends crucially upon each of these items. Identification of features refers to discovering the relevant characteristics for discriminating between the classes. For example features may include morphological characteristics such as size (area, perimeter), shape (ellipticity, circularity, moments), optical characteristics such as intensity distribution within the object, and colour, and transform domain characteristics using discrete Fourier transforms, discrete cosine transforms, wavelet transforms, eigen-vector based transforms or the Karhunen-Loève transforms. Choosing features with good discriminating ability is essential to the success of the classifier. Over the last 30 years, significant progress has been made in the theory and design of pattern classifiers. The re-emergence of artificial neural network techniques has contributed to several new techniques towards the design of classifiers. Prior to this, the dominant paradigms in use were statistical based, structural or the syntactic based classifiers and data based such as the k-nearest neighbour classifiers.

In spite of all these advances, very little progress has been made towards evaluating different classifiers. Often, known theoretical results on error bounds and probability of errors based on normal distribution of class conditional densities and/or infinite samples are used in the performance evaluation and the design. A pure theoretical evaluation of classifiers for finite samples is very difficult. Moreover, the performance of a

classifier may vary from one set of data to the other, and further one classifier may outperform the others for one set of data while it may perform poorly for the other set of data. Therefore for a finite sample, one should consider a suite of classifiers rather than a single classifier both for design and implementation.

In this paper, design, evaluation and implementation of a suite of classifiers is addressed when a training set of finite samples is available. In the literature a number of pattern classification schemes have been proposed which may be broadly classified into k- nearest neighbour classifiers, statistical pattern classifiers and artificial neural network classifiers. All of these have their strengths as well as their weaknesses [1-7].

Depending upon the application, one scheme may prove to be better than other. A classifier which gives the *best performance* is selected from a number of state of the art pattern classification schemes which are diverse in the design approach as well as the assumptions employed in their design. The selected classifiers include the KNNC, the MMDC and the ANNC. The performance of all these classifiers approach the Bayes error rate as the sample size becomes infinitely large. The main hurdle in evaluating the performance of a classifier is that an explicit analytical expression for the Bayes error rate is too difficult to find even for an infinite data set, and even if it is found, it may be too complicated for design purposes; hence instead of the Bayes error rate, its bound is employed for evaluating the performance of a classifier, and the bound is obtained from a measure of distance between the two probabilities density functions [6]. A number of measures of distance have been proposed in the literature including the Bhattacharyya distance measure, the Mahalanobis distance measure, and the Kullback-Liebler divergence measure [3].

Depending upon the data set, one measure may give a tighter bound than the other. A bound, which is the Bhattacharyya distance measure when the underlying pdfs are Gaussian, is used to measure the performance of the classifier. This bound is termed herein as the *benchmark*. From extensive simulation results, it was found that when the underlying pdfs are uni-modal and symmetrical, the benchmark gives a fairly tight bound on the Bayes error rate. A classifier with the best performance is selected from a list of classifiers as follows:

- The classifiers are grouped in the order of their increasing classification performance using the bounds which relate the classifier performance with that of the Bayes error rate. For example, the ANN is placed at the top and the KNN or MMDC is placed at the bottom.
- The classifier at the bottom of the hierarchy is designed first using the benchmark as an achievable target. The next classifier in the hierarchy is designed by using the performance of the previous classifier as well as the benchmark as achievable targets, and so on till the classifier at the top of the hierarchy is designed. The hierarchical scheme will ensure that each classifier is 'optimal'.

HIERARCHICAL APPROACH TO CLASSIFIER DESIGN

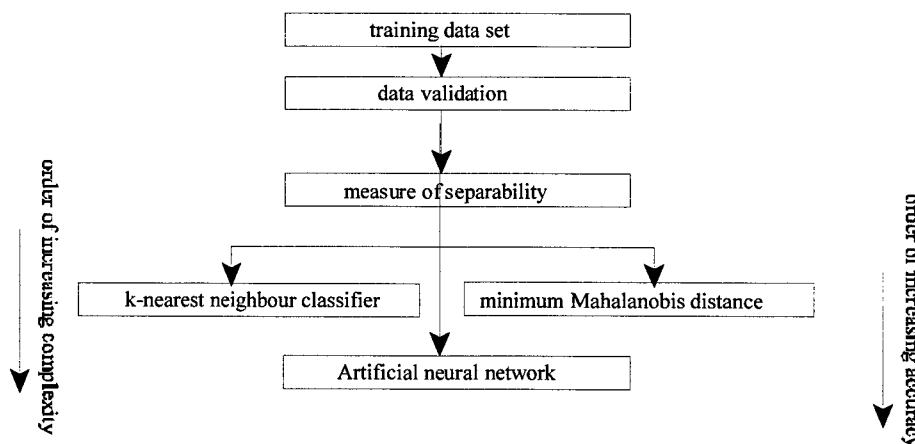


Fig. 2. The hierarchical structure: a methodology for the design of a pattern classifier.

The classifier which is reliable and yields a superior performance is obtained by adopting a hierarchical approach. The pattern classification task is divided into a number of tasks and these tasks are executed in a well defined order. The tasks are Data Analysis Scheme (data validation, data visualization and feature quality analysis), Pattern Classification Schemes (the KNNC, the MMDC, and the ANNC). The classification scheme should be properly designed and correctly implemented for the given application. In the design of the KNNC, the choice of the number of training data elements, K , closest to a test data element should be selected so as to achieve a high classification rate without increasing the computational burden. In the MMDC, the mean and covariance of the data must be accurately estimated. In designing an ANN, the architecture of the ANN must be tailored to the given problem. In the case of unlimited training data, all the above classifications schemes will approach the lowest achievable error rate sometimes referred to as the *Bayes error rate*. Using the well known measure of class separability, namely the Bhattacharyya distance, an upper and lower bound on the error rate are computed from the training data set. This distance measure is used as a bench mark to evaluate the performance. As the derived bounds may not be tight, it is difficult to gauge whether or not the classification scheme has attained its peak performance.

An additional check is provided by cross-checking the relative performances of the various pattern classification schemes. Generally, the higher the design complexity, the more accurate is the classification scheme. The pattern classification schemes are executed in the order of increasing accuracy and increasing complexity so that during the design phase, a task at one level sets the performance goal (e.g. classification accuracy) for the task at the next level: the KNNC and the MMDC require less design effort than the ANNC and thus their performance will be used as achievable target for the performance of the ANNC. The KNNC is the least complex in its design as the entire data set rather than parameters estimated from the data set are used. Further the KNNC is better than the MMDC if the underlying pdf is asymmetrical and/or multi-modal and particularly when the class separation is crisp, while the MMDC is better than the KNNC if the underlying pdf is unimodal and symmetrical. This approach of cross-checking the performances will ensure that the design and the implementation of the classifiers are appropriate for the problem. Fig. 2 illustrates the hierarchical approach to the design of the proposed scheme.

DATA ANALYSIS

The proposed data analysis scheme consists of the following: data validation, data visualization, feature quality analysis. The reliability of a pattern classification scheme will be no better than the reliability of the data used to estimate the parameters of the scheme (e.g. the parameters describing the separating surfaces of the MMDC or the weights of the ANN). The common source of errors in the training set includes mislabelling of the data (the data may be classified erroneously as belonging to a class C_i when in fact it belongs a different class C_j), malfunctioning of the data acquisition system, and errors in the numerical computation of the features. The goal of data validation is to flag any suspicious feature values. Spotting any erroneous values before proceeding to develop a classifier saves an enormous amount of time. Data validation essentially involves the computation of the mean μ_{ij} and variance σ_{ij} of each class C_j . A data point which falls more than 3 standard deviations from the class mean is flagged as an outlier or a bad data: a data point x_i is an outlier if

$$|(x_i - \mu_{ij})| > 3\sigma_{ij}, \quad x_i \in C_j, \quad \sigma_{ij}$$

When an outlier is detected, the data is cross-checked to ensure that the data is not mis-labelled, the data acquisition system is not faulty or there are no errors in computing the features. Note that the outlier is not discarded as it might not be a bad data point. Data visualization is becoming an increasingly important tool for understanding, interpreting, analysing, and validating the data. In the present context, it provides a powerful qualitative picture of the separability of given data into various classes. A simple scheme to visualize the class separability is to orthogonalize the given set of features. From the mean and the covariance of the data, the given set of features was orthogonalized using the Singular Value Decomposition (SVD). The whitened features are uncorrelated. The whitened features are plotted to provide a visual measure of class separability and thus provide a feel for the complexity of the classification task. An intuitive measure of class separability is the difference in the means between any

two classes compared to their covariances: the larger the difference in the mean compared to their covariances the better the class separability. This intuitive concept of class separability is formalized by statisticians to derive a number of measures such as the Kullback-Liebler divergence, the Bhattacharyya and the Mahalanobis distances[2-6]. In this work, the Bhattacharyya distance measure is employed. The Bhattacharyya distance measure, denoted $B(i,j)$, is a distance between two conditional pdfs $f(x|H_i)$ and $f(x|H_j)$ where H_i and H_j are the hypotheses, and is given by

$$B(i,j) = -\log \rho, \text{ where } \rho = \int_{-\infty}^{\infty} \sqrt{pf(x|H_i)(1-p)f(x|H_j)} dx$$

Since ρ lies between 0 (when the two pdfs are non-overlapping) and $\sqrt{p(1-p)}$ (when the two pdfs are identical). The Bhattacharyya measures satisfy all the postulates for a distance measure or metric except the triangular inequality. For the Gaussian case $B(i,j)$ reduces to

$$B(i,j) = \frac{1}{8} (\mu_i - \mu_j)^T \Sigma^{-1} (\mu_i - \mu_j) + \frac{1}{2} \log \left(\frac{|\Sigma|}{\sqrt{|\Sigma_j||\Sigma_i|}} \right) - \frac{1}{2} \log(p(1-p))$$

where $\Sigma = p \Sigma_i + (1-p) \Sigma_j$. The Bhattacharyya distance is a sum of two terms. The first term gives the class separability due to the mean difference while the second term gives class separability due to the differences in the covariances. The first term vanishes when $\mu_i = \mu_j$, while the second term vanishes when $\Sigma_i = \Sigma_j$. For the binary case, the Bayes error rate R^* is related to the Bhattacharyya distance measure by the following inequality given by [6]:

$$p - \frac{1}{2} \sqrt{1 - 4\rho^2} \leq R^* \leq p - \frac{1}{2} + \rho \quad \text{for all pdfs.}$$

All these inequalities are tight; there exists a pdf for which the upper or the lower bound given by the above inequalities is attained.

PATTERN CLASSIFICATION SCHEMES

The Bayes classification scheme cannot be employed as the pdf is not generally known *a priori*, or due to limited data it is not possible to estimate accurately the pdf. One has to settle for non-optimal classification schemes which are designed based either on the assumed pdf or on the assumed form of the data clusters (which is related to the underlying pdf governing the data). The MMDC is based on the assumption that the pdf is Gaussian while the KNNC assumes that the data points form clusters containing mostly samples from the same class. The MMDC is indicated even when the pdf is *thick-tailed* and non-Gaussian as long as it is uni-modal and symmetrical, while the KNNC can handle both multi-modal and non-symmetrical pdf as long as it is *thin-tailed*. The ANNC is based on an entirely different approach. The ANNC will be able to classify any arbitrary set of patterns, with minimal *a priori* information about the data. However, the ANN architecture must be sufficiently complex, and the training set size must be sufficiently large and must be representative of the population to which it will be ultimately applied. All the classifiers namely the KNNC, the MMDC and the ANNC have their strengths as well as their weaknesses. Any one of the classifiers may outperform the rest depending upon the problem. Hence one has to analyse the performance of all the classifiers to select an appropriate one. See Table 1.

EVALUATION OF THE PROPOSED SCHEME

The proposed systematic and reliable approach was evaluated on the problem of classification of biological cells. The objective of the work was to develop a two-way classification to assert a given cell as NORMAL or ABNORMAL. The data set was randomly divided into two equal-size training and test sets: the NORMAL containing 1621 cells and the ABNORMAL containing 1117 samples. The table gives the performance of the various classifiers. The given pre-classified data was subject to data validation, data visualization and feature quality analysis. There were 10% outliers and the data visualization did not yield a clear picture of class separation. The outliers were cross-checked, and it was found that they occurred mainly as a result of mis-labelling in the training set. The mis-labelled data was corrected, and correctly labelled training data set was employed in this work. The overall Bhattacharyya distance separating the

classes namely NORMALS and the ABNORMALS was computed to be 1.3 when the underlying pdfs are assumed to be Gaussian. This corresponds to an overall classification rate of 86% , and this was used as the benchmark. The upper bound of the classification is $(86+4)\%=90\%$ and the lower bound on the calcification accuracy is $(86-4)\%=82\%$. Thus the estimates of the minimum, $R_{\min} = (100 - 90)\% = 10\%$, the maximum $R_{\max} = (100-82) \% = 18\%$, and the benchmark, $R_0 = (100-86)\% = 14\%$. A validity check for R_{\min} and R_{\max} is given by the following inequality

$$R^* \leq \frac{2p(1-p)}{1+p(1-p)\Delta^2(i,j)} \Rightarrow R_{\min} \leq 22\% \text{ where } \Delta^2(i,j) = (\mu_i - \mu_j)^T \Sigma^{-1} (\mu_i - \mu_j)$$

Table.1. Comparison of classifiers: binary classification

Classifier	probability of error	Comment
Bayes	$p - \frac{1}{2}\sqrt{1-4p^2} \leq R^* \leq p - \frac{1}{2} + p$ $R^* \leq \frac{2p(1-p)}{1+p(1-p)\Delta^2}$	The Bayes error rate is the minimum achievable error rate. Its upper and lower bounds are functions of the Bhattacharyya distance measure. An additional upper bound is available which is a function of the Mahalanobis distance.
KNNC	$R_{NN} \leq 2R^*(1-R^*)$	The error rate, R_{NN} , is no more than twice the Bayes error rate. The data is assumed to form dense clusters containing mostly samples from the same class. It can handle multi-modal and asymmetrical pdfs as long as they are thin-tailed.
MMDC	$R^* \leq R_{MMDC}$	The error rate, $R_{MMDC} = R^*$ if the pdfs are Gaussian with identical covariances. It can handle thick-tailed non-Gaussian pdfs as long as they are uni-modal and symmetrical.
ANN	$R^* \leq R_{ANNC}$	If the training set is sufficiently rich and the ANN architecture is sufficiently complex, and the training algorithm is successful, the error rate $R_{ANNC} = R^*$.

The data analysis was found to be very crucial to the classifier design. Many errors due to mis-labelling of the pre-classified data and numerical errors in the computation of the feature were detected before the data was used in the classifier design. The KNNC is the simplest to design as it assumes no a priori information and no statistical parameters are computed from the data. However, the KNNC suffers from the computational complexity: larger the k value the more complex the algorithm. There exists a trade-off between the complexity and the accuracy in the choice of k. With $k=3$ the error rate was found to be 15.7%. The KNNC is designed first to serve as a performance target for the other classifiers. The MMDC classifier requires the estimation of the mean and the covariances. In the computation of the covariances, the Bhattacharyya measure was used to decide between the hypothesis that the covariances are equal against the hypothesis that they are different.

Since

$$\frac{|\Sigma|}{\sqrt{|\Sigma_j||\Sigma_i|}} \approx 1$$

the covariances were assumed to identical and the MMDC classifier was reduced to the *Linear Discriminant Classifier*. The design error of the MMDC compared to that of the KNNC stem from the errors in the estimation of the mean and more so from the covariances. Hence the error rate of the KNNC is used to cross-check the performance of the MMDC. Then the design of the ANNC was considered. The error rate of the MMDC was found to be 15.7%. The ANN is limited to two-layers with the number of nodes in the hidden layer equal to 6. The benchmark $R_0 = 14\%$, the error rate of the KNNC, 15.7% and the error rate of the MMDC, 15.7% served to chose an appropriate architecture for the ANNC as well as to verify whether the network has been properly trained. The appropriateness of the design of the classification scheme and the correctness of its implementation were verified. The reliability of the design

and the implementation of the classifiers were verified by ensuring that the accuracy of the classifiers satisfy the following *reliability criteria*

- I. $\text{acc}_1 \approx \text{acc}_2 \leq \text{acc}_3$
- II. $82 \leq \text{acc}_i \leq 90, i=1,2,3.$

where acc_1 is the accuracy of KNNC, acc_2 is the accuracy of the MMDC, acc_3 is the accuracy of the ANNC. The first reliability criterion is derived from the theoretical justification and the assumptions used in the design of the classifiers while the second criterion is based on the upper and the lower bounds on the Bhattacharyya distance. It should be noted that the Bhattacharyya distance measure predicts a classification accuracy of $86 \pm 4\%$. All the results have an associated accuracy in the range which ensures the design and implementation of the three basic classification algorithms is appropriate and reliable. Also from the results it is seen that the accuracy of the KNNC is comparable to that of the MMDC whereas the ANNC outperforms the other two classifiers. Table 2. gives the performance of the various classifiers.

Table 2. performance of the various classifiers

CLASSIFIER	error rate	Accuracy	validity
1. The k-nearest neighbour classifier (KNNC)	15.7%	84.3%	$R_{NN} \leq 2R_{\min} (1-R_{\min}) : 15.7\% \leq 18\%$ $R_{NN} \leq R_{\max} : 15.7\% \leq 18\%$
2. The minimum Mahalanobis distance classifier (MMDC)	15.7 %	84.3 %	$R_{\min} \leq R : 10\% \leq 15.7\%$ $R \leq R_{\max} : 15.7\% \leq 18\%$
5. The artificial neural network classifier(ANNC)	13.08%	86.92 %	$R_{\min} \leq R : 10\% \leq 13.08\%$ $R \leq R_{\max} : 13.08\% \leq 18\%$

CONCLUSIONS

A hierarchical approach to the selection of an appropriate classifier which is reliable and yields a very high accuracy is proposed. The results of its evaluation based on actual data is highly encouraging. The classifiers were chosen to be the KNNC, the MMDC and the ANNC. They are designed using entirely different methodologies. This will serve to eliminate common mode misclassification errors. The data validation scheme was able to flag suspicious data before designing and implementing a classifier. The benchmark, an estimate of the minimum and maximum achievable error rate served to validate the design of the classifiers. The reliability of each of the classification schemes was verified by ensuring that the classifier at a lower level in the hierarchy has a higher classification accuracy and the accuracy of each of the classifiers is within upper and the lower bounds.

ACKNOWLEDGEMENT

The authors acknowledge the support of Natural Sciences and Engineering Research Council of Canada. The authors appreciate the help of Mr. Ram Balasubramanian.

REFERENCES

1. R.O. Duda, P.E. Hart, 1973. Pattern Classification and Scene Analysis, John Wiley and Sons.
2. K. Fukunaga, 1990. Introduction to Statistical Pattern Recognition, Academic Press, San Diego.
3. Luc Devroye, Laszlo Györfi, Gabor Lugosi, 1996. A Probabilistic Theory of Pattern Recognition, Springer Verlag, NY.
4. Robert Schalkoff, 1992. Pattern Recognition: statistical, structural and neural approaches, John Wiley and Sons.
5. Michele Basseville, Igor V. Nikiforov, 1993. Detection of abrupt changes: theory and applications Prentice-Hall, NJ.

6. Thomas Kailath, 1967. The divergence and Bhattacharyya distance measure in signal detection. IEEE Transactions in Communication Technology, 15, 52-60.
7. T.M. Cover, P.E. Hart, 1967. Nearest Neighbour classification, IEEE Transactions on Information Theory, IT-13(1), 21-27.

Dynamic Associative Memory using Chaotic Neural Networks

Y. Fukuhara, Y. Takefuji

Keio University, Graduate School of Media and Governance,
5322 Endo, Fujisawa, 252-0816, Kanagawa, Japan

Tel: 81 0466 49 1062 Fax: 81 0466 47 5125 Email: fuku@mag.keio.ac.jp

ABSTRACT

In this paper, we propose a multi-module chaotic associative memory (MCAM) that uses chaotic neural networks. In this method, the chaotic associative memories are connected to each other. If MCAM can not obtain enough information of a target, MCAM shows a behavior that looks like human "perplexity", where MCAM succeeds in one-to-many associations. And when MCAM obtains enough information to recognize a target, MCAM converges to a stable state. Although the structure of MCAM is simple, MCAM realizes one-to-many association by using chaotic dynamics.

INTRODUCTION

The purpose of this research is to simulate brain-like information processing using chaotic neural networks. An object contains a variety of information such as shape, color, smell, etc. Humans can recognize an object by obtaining partial information associated with the object. Suppose you see someone waving his hand at you but you cannot see his face clearly. In that situation, you try to identify the person according to your memory based on only partial information you receive, for example, the shape, height, clothes, or gestures of the person. It is only when you find a definitive clue identifying the person that you will be able to recognize him with confidence. When he calls you, you will further try to guess who he is by his voice and body shape, then identify him with the voice information newly obtained. When we can obtain only partial information of the target, we associate many things that are related to that information. If we have enough information to recognize someone, we narrow down the search domain of the selection.

We use a chaotic neuron model. Actually, many physiologists report that chaotic dynamics are observed in a biological neuron [8]. Therefore, we think the feature of chaotic dynamics is significant for artificial neural networks. Chaotic associative memory models, related to the proposed system, have been developed. Conventional models converge to certain patterns and wander to other patterns one after another if only a single chaotic associative memory is used. Generally, control of chaotic networks is an intractable problem.

To control chaotic behavior, we combine several chaotic associative memories together. Each associative memory is assigned respectively to one aspect of information such as shape, voice, smell, etc. Due to the chaotic dynamic capability, MCAM can associate one-to-many relations, if the system cannot obtain enough information, and then MCAM can obtain additional information that is enough to choose the target among candidates. The system can converge to one state immediately. This stable state means that "the entire information of an object is synchronized".

The most significant advantage of our system is that we can control the chaotic behavior intelligently. MCAM performs well under a noisy environment in practical use, even when the system cannot obtain enough information of the target.

METHOD

In this section, we introduce a chaos neural model, the chaotic associative memory, and the structure of MCAM. We explain how to manipulate MCAM.

Chaos-Neural Networks

K.Aihara proposed a chaotic neural model [1, 2]. The model imitated a chaotic behavior observed in a biological neuron, such as refractoriness. Dynamics of chaotic neural model is given by:

$$x_i(t+1) = f\left(\sum_{j=1}^M V_{ij} \sum_{d=0}^t K_e^d A_j(t-d) + \sum_{j=1}^N W_{ij} \sum_{d=0}^t K_f^d x_j(t-d) - \alpha \sum_{d=0}^t K_r^d x_i(t-d) - \theta_i\right) \quad 1.$$

$$f(u) = 1/(1 + \exp(-u/\epsilon)) \quad 2.$$

where $x_i(t+1)$ is the output of a chaotic neuron at time $t+1$, M is a the number of given data, N is the number of neurons, $A_j(t)$ is j th input value, V_{ij} is a weight from A_j to i th neuron, W_{ij} is a weight from i th neuron to j th neuron, and θ is a threshold. Note that α, K_e, K_f, K_r are constant numbers. Eq. 2. is a sigmoid function and ϵ defines the slope of Eq. 1. If all input values A_j are kept as one state, Eq. 1. can be translated into the following three formulas:

$$xi(t+1) = f(\eta i(t+1) + \zeta i(t+1)) \quad 3.$$

$$\eta i(t+1) = k\eta i(t) + \sum_{j=1}^N w_{ij}x_j(t) \quad 4.$$

$$\zeta i(t+1) = kr\zeta i(t) - \alpha xi(t) + ai \quad 5.$$

where ζ is a term of the input value and η is a term of the mutual interactions. Note that ai is a constant number. Also proposed was an associative memory combining several chaotic neurons similar to the Auto-Associative Memories[3]. The synapse weight W is given by:

$$w_{ij} = \frac{1}{Q} \sum_{p=1}^Q (2x_i^p - 1)(2x_j^p - 1) \quad 6.$$

where Q is the number of patterns. The chaotic associative memory can memorize some patterns by using Eq. 6. Due to the effect of chaotic dynamics in the neurons, a chaotic associative memory wanders between several states. To simulate this brain-like behavior, we use this chaotic associative memory as part of a MCAM.

System Structure of MCAM

In this section, we explain how to connect several chaotic associative memories. This is the most important aspect of this paper. Figure 1 shows the system structure of MCAM. In our system, several chaotic associative memories are combined together. This structure is similar to the Multidirectional Associative

Memory (MAM)[4]. Each of the memories can accept information from the outside and from their neighboring memories at all times.

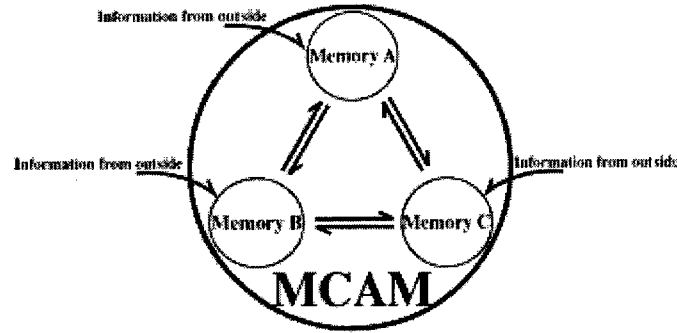


Fig. 1. Structure of MCAM

Based on Eq. 4., η_i is given by:

$$\eta_i(t+1) = \begin{cases} k\eta_i(t) + \sum_{j=1}^N w_{ij}x_{mj}(t) & \left(\sum_{j=1}^N I_j = 0 \right) \end{cases} \quad 7a.$$

$$\eta_i(t+1) = \begin{cases} k\eta_i(t) + \sum_{j=1}^N w_{ij}(\tau x_{mj}(t) + (1-\tau)I_j) & \left(\sum_{j=1}^N I_j \neq 0 \right) \end{cases} \quad 7b.$$

$$\eta_i(t+1) = \begin{cases} k\eta_i(t) + \sum_{j=1}^N w_{ij}(\tau x_{mj}(t) + 1/g(1-\tau)(\sum_{m' \in M-m} \phi(x_{m'j}(t)))) & \left(\sum_{j=1}^N I_j = 0 \text{ and } \sum_{j=1}^N |x_{m'j}(t) - P_{m'l}| = 0 \right) \end{cases} \quad 7c.$$

$$\eta_i(t+1) = \begin{cases} k\eta_i(t) + \sum_{j=1}^N w_{ij}(\tau x_{mj}(t) + (1-\tau)(\kappa I_j + 1/g(1-\kappa)(\sum_{m' \in M-m} \phi(x_{m'j}(t)))) & \left(\sum_{j=1}^N I_j \neq 0 \text{ and } \sum_{j=1}^N |x_{m'j}(t) - P_{m'l}| = 0 \right) \end{cases} \quad 7d.$$

$$x_{m'j}(t) = \begin{cases} 1 & (x_{m'j}(t) > 0) \\ 0 & (x_{m'j}(t) \leq 0) \end{cases} \quad 8.$$

$$\phi(x_{m'j}(t)) = \left(\sum_{s \in Q} P_{msj} \right) / n \quad \left(\sum_{j=1}^N |x_{m'j}(t) - P_{m'l}| = 0 \text{ and } P_{m'l} \text{ is related with } P_{ms} \right) \quad 9.$$

where $x_{mj}(t)$ is an output of the j th neuron in the m th memory, I_j is the information from outside, N is the number of neurons, M is the set of all memories, m' is a set of the memories associated with any pattern, g is the number of memories that are associated with any pattern, $x_{m'}(t)$ is a vector of the m' th memory and $P_{m'l}$ is a vector of the l th pattern that is memorized in the m' th memory, S is a set of related patterns, Q is a set of memorized patterns, P_{msj} is the j th information in the pattern P_{ms} which is memorized in the m th memory, and n is the number of related patterns. Note that τ and κ are constant numbers.

Eq. 7. is divided into the following four conditions:

1. If the memory obtains no information from outside (Eq. 7a.)
2. If the memory obtains any information from outside (Eq. 7b.)

3. If the memory obtains no information from outside and other memories associate something (Eq. 7c.)
4. If the memory obtains any information from outside and other memories associate something (Eq. 7d.)

Defining Relations between the Memories

We have to define relations between the patterns beforehand. There are two kinds of relations: one-to-one and many-to-many. When a memory converges to a certain pattern, other memories obtain information dependent on the relation (Eq. 9.).

One-to-One Relations

If a pattern $a1$ in a Memory A is related to a pattern $b2$ in Memory B, then as Memory A converges to pattern $a1$, pattern $b2$ is given to Memory B.

One-to-Many Relations

If the pattern $a1$ which is memorized in Memory A is related to patterns $b2$ and $b3$ contained in Memory B, the arithmetic mean between $b2$ and $b3$ is given to Memory B.

SIMULATION AND RESULTS

To show the effectiveness of MCAM, we experimented with a computer simulation. In this simulation, MCAM contains three chaotic associative memories with each containing one hundred chaotic neurons and a memory of three patterns. MCAM learned the relationships and used the parameters shown in Figure 2.

1. pattern $a2$ relates to pattern $b3$ and $c3$, i.e., a one-to-one relation.
2. pattern $a1$ relates to pattern $b1$ and $c1$, and also, to $b2$ and $c2$, i.e., a one-to-many relation.

We show three experimental results for these features of MCAM.

Simulation 1: One-to-Many Associations

First, only pattern $a1$ is given to Memory A as initial data, while Memory B and Memory C are empty as shown in Figure 3.

At state (I), we observe that Memory A recalls pattern $a1$, Memory B recalls the pattern $b2$ and, Memory C recalls pattern $c2$ and $c1$ by changing the state in turn. At state (II), Memory B recalls pattern $b1$ and Memory C recalls pattern $c1$ synchronously. At state (III), Memory B recalls pattern $b2$ and Memory C recalls pattern $c2$ synchronously. And at state (IV), Memory B recalls pattern $b1$ and Memory C recalls pattern $c1$ similar to state (II).

As a result of this simulation, we confirm that MCAM can succeed in one-to-many association, even when the system does not have enough information.

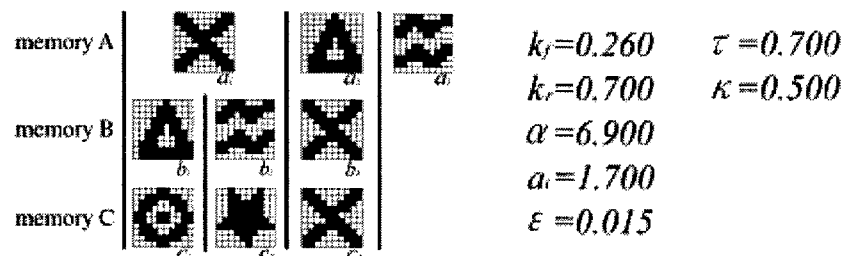


Fig. 2. Relationships between the memorized patterns and parameters in MCAM

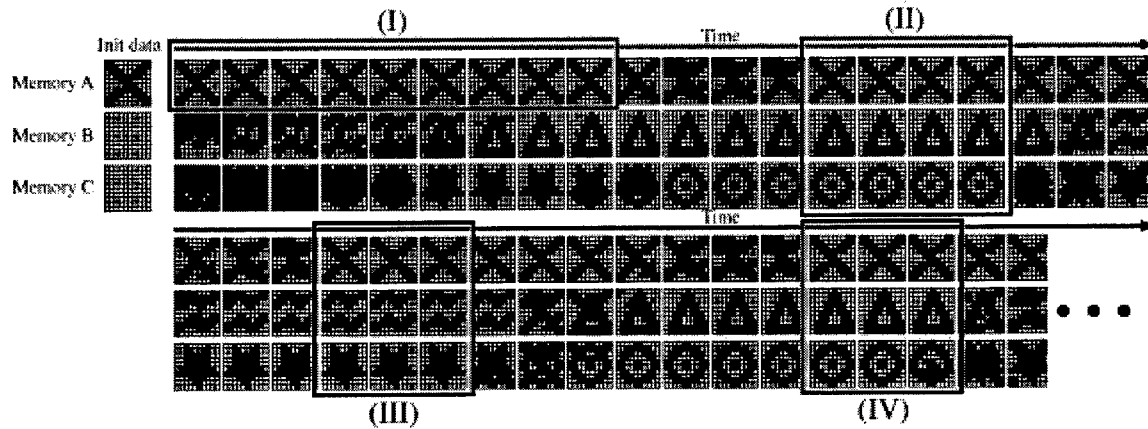


Fig. 3. One-to-many relation.

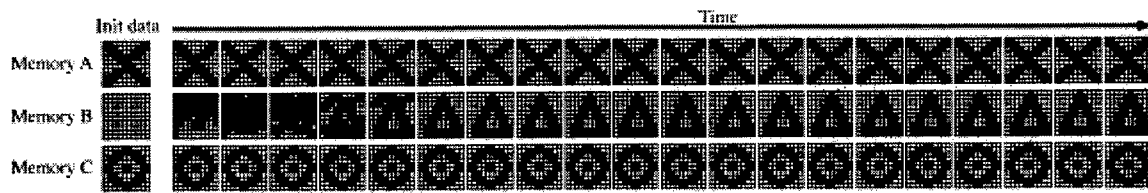


Fig. 4. Complementary Information

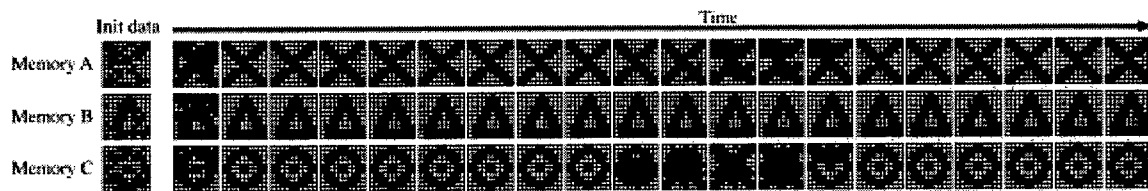


Fig. 5. Input with noise.

Simulation 2: Complementary Information

Pattern a1 is given to Memory A and pattern c2 is given to Memory C as initial data, while Memory B is empty, as shown in Figure 4. We can observe that Memory A recalls pattern a1, Memory C recalls pattern c2 and Memory B recalls pattern b2, and the system keeps the state stable. As a result of this simulation, if MCAM has enough information, MCAM can converge to a single stable state.

Simulation 3: Inputs with Noise

The inputs with noise are given to the memories, as shown in Figure 5, and MCAM tries to converge to a stable state. It depends on the quantity of noise whether or not the memories converge to the correct pattern. If the given data contains large noise, the memories converge to irrelevant patterns once in a while, but in most case, the memories converge to correct patterns. As a result of this simulation, MCAM performs well even if the given inputs are relatively noisy.

CONCLUSION

The multi-module chaotic associative memory (MCAM) can provide one-to-many associations using chaotic dynamics. CBAM[5], CMAM[6] and MMA[7] are related to our models. Although CBAM uses chaotic neurons to express a context information, MCAM does not need that information. MMA is superior

to our model regarding one-to-many or many-to-many association, but the structure of MCAM is easier and more natural than MMA.

If MCAM cannot obtain the information necessary to converge to a stable state, MCAM shows an interesting behavior like "perplexity", i.e., the system wanders until new information is fed into the model. When new information is input, the system can immediately and constantly utilize the new information, while conventional systems must compute from the initial state again. In addition, MCAM performs well under noisy environment. The simulation results show that MCAM succeeds in the one-to-many association. In the future, MCAM will be able to provide many-to-many associations, if the chaotic associative memory can memorize a large number of patterns.

REFERENCES

1. K.Aihara, T.Takabe and M.Toyoda, 1990. Chaotic Neural Networks, *Phys. Lett. A*, 144, 6, 7, 333-340
2. K.Aihara, 1990. Chaotic Neural Networks, in "Bifurcation Phenomena in Nonlinear Systems and Theory of Dynamical Systems" (H.Kawakami ed.), 143, *World Scientific*, Singapore
3. J.J.Hopfield, 1982. Neural Networks and Physical Systems with Emergent Collective Computational Abilities, *Proc. Of National Academy of Science, U.S.A.*, 79, 2445-2558
4. M.Hagiwara, 1990. Multidirectional Associative Memory, *IJCNN, Washington, D.C.*, 1, 3
5. Y.Osana, M.Hattori and M.Hagiwara, 1996. Chaotic Bidirectional Associative Memory, *ICNN, Washington D.C.*, 2, pp.816-821
6. Y.Osana, M.Hattori and M.Hagiwara, 1997. Chaotic Multidirectional Associative Memory, *ICNN, Houston*, 2, pp.1210-1215
7. M.Hattori and M.Hagiwara, 1998 Multimodule associative memory for many-to-many associations, *Neurocomputing* 19, 99-119
8. C.A.Skarda and W.J.Freeman, 1987. How brains make chaos in order to make sense of the world. *Behav, Brain Sci.* 10, 161-195.

Trends in Intelligent Process Control Methods in the Primary Aluminum Industry

R.T. Bui*, L.G. Tikasz*, J. Perron**

* Université du Québec à Chicoutimi, Chicoutimi, Quebec, Canada, G7H 2B1

** Alcan International Limited, Jonquière, Quebec, Canada, G7S 4K8

ABSTRACT

This progress report presents new trends and work underway in research related to process control methods for the primary aluminum industry, based on process modeling combined with advanced control techniques using computational intelligence. Promising applications are seen for a wide range of process control situations from calcining kilns to electrolytic cells and casting furnaces.

INTRODUCTION

With the advent of computational intelligence (CI), coupled with the now well accepted model-based control, new trends are opening up for process control in the aluminum industry. Mathematical models run on computers serve as process simulators providing a convenient, low cost, user-friendly and risk-free alternative to the traditional trial-and-error performed on real processes. Models are nowadays present in laboratories as well as on plant floors, and are used as tools for research, training, process analysis, parameter studies, and even process supervision. On the other hand control emulators, also run on computers, emulate the process's control system, providing a tool for analyzing and evaluating new control schemes. Depending on the process, control schemes may apply off-line or on-line, open-loop or closed-loop techniques, involving human intervention to various degrees, and based on conventional procedures of the PID type, or more recently and still in a rather primitive way, on one form or another of the new techniques loosely identified as Artificial Intelligence (AI) which notably includes knowledge bases, expert systems, neural networks, fuzzy logics or genetic algorithms. In some cases, a combination of more than one categories is required to do the work.

The primary aluminum industry is endowed with a remarkable wealth in thermo-physical processes which require a wide range of process control techniques. Along the protracted progression leading from bauxite to aluminum alloys, a variety of processes are needed, each of which assorted with tight criterias on quality and giving rise to high added values for the end products. We are in fact dealing with a chain of transformation processes, the end product of one process serving as incoming material for the next process down the line. The added values benefiting the end product in some cases amount to several times the cost of the incoming material. Suffice it to think of examples such as commercial alumina resulting from the processing of trihydrates, or anode blocks resulting from the processing of green petroleum coke, or cathode blocks from green anthracite, or metal matrix composites from aluminum.

Good control not only ensures the productivity of the process and the quality of the product, it also yields a better utilization of the high-capital equipment and lengthens its life duration. This is particularly evident for processes operating in hostile environments in terms of corrosion and high temperature, such as the aluminum electrolytic cell. A good control will also have positive impact on the ecology and the workers' quality of life by reducing the harmful discharges into the environment. A few cases in point with easily measurable benefits can be found in the reduction of carbon dust discharge into the atmosphere by the petroleum coke calcining kilns, or reduction of the volume of the red mud rejected by the Bayer process that extracts alumina trihydrates from bauxite. Those are indeed convincing reasons for the industry to focus effort on process control methods, especially at this particular juncture in time when the ubiquitous computer offers innovative solutions for communication, information and networking.

STATE OF THE ART

Most control systems presently used in the aluminum industry apply the conventional techniques with PID or PD feedback either in open loop or closed loop. The resulting control algorithms are generally programmed into computers or microcomputers and fine-tuned on the real process. In some cases, PID control is combined with fuzzy logics. A current example can be found in a 1996 U.S. Patent [1] proposing a control scheme for the rotary coke calcining kiln. The system measures the temperatures at various positions along the kiln and adjusts the control variables to move the calcining zone to the desired optimal position. As part of the system, a fuzzy logic controller determines the flowrate of complementary air and the rotational speed of the kiln. In some other cases, conventional control techniques are applied in combination with a knowledge base, which amounts to a simple form of expert system in which the many elements of the heuristic and fundamental knowledge coming from scientists, process designers, engineers and operators, cumulated over the years, are organized, structured, and encapsulated by knowledge engineers into computer programs in the form of rules and decision trees to serve as off-line consultation for the operators [2]. Note that this involves an intimate knowledge of the process, and more often than not the relevant details are regarded as highly proprietary; as a result it is unlikely that such a product would give rise to publications in the open literature.

Process control systems that involve some aspects of the so-called advanced techniques are often accepted reluctantly or even rejected by plant floor workers. This is understandable as operators, working in already stressful conditions, prefer a familiar environment. This human element must be dealt with as an inherent part of the problem. In fact, the Aluminum Technology Roadmap Workshop of 1996 explicitly recommended research on process control as one of the priorities for the years to come [3].

The other major aspect of process control involves the modeling of processes (Figure 1). As a safe, convenient and low-cost alternate to trials-and-errors performed on the real equipment, mathematical models operated on computers are fast becoming the rule in process industries as tools for process control and supervision as well as for training. Depending on the nature of the process and the purpose of modeling, a variety of models can be found, from lumped-parameter, steady-state models to distributed-parameter, dynamic, multidimensional models, the latter are often complex enough to require several man-years of hard work to build, sophisticated CFD numerical codes to solve, and powerful computers to run. What is more, mathematical modeling often needs to be complemented with physical modeling, in the form of laboratory size prototypes or small-scale models, to help tackle those aspects of modeling that still elude a formal understanding either in rigorous form or in empirical form.

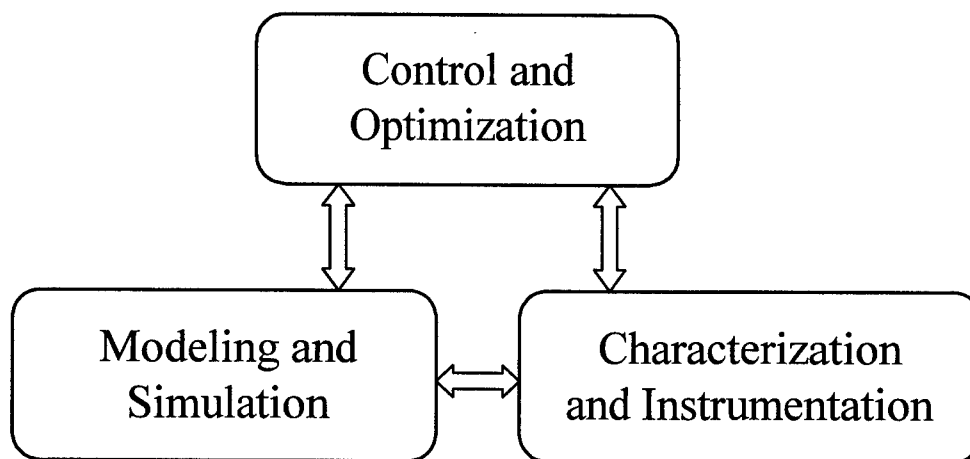


Fig. 1. The three facets of process control.

Mathematical models require inputs in order to run and produce the desired outputs. Inputs are in the form of materials properties such as thermal or electrical conductivities, specific heats, densities, or in the form of parameters qualifying the transport mechanisms such as heat transfer coefficients, emissivities, reaction rates or viscosities. This leads to another crucial facet of process control, commonly referred to as characterization (Figure 1). Clearly, it makes sense to attempt to control a process only if it is possible to characterize it. This calls for the development of measurement techniques, a field of research of considerable importance in its own right. In quite a few cases, instruments (sensors) are developed to measure the variables directly; in other cases, direct measurements are not possible, and the needed variables can only be evaluated through calculations based on other variables that can be measured. These "virtual sensors" have become a challenging branch of activity of utmost interest to process engineers.

The above three research domains namely process control, process modeling and process characterization, have been at the center of research activities for quite some time. The new element is that with the advent of new information and communication technologies, the above three domains can be put to work in a complementary manner, enabling researchers, designers, operators and managers to collaborate in real time, ignoring geographical barriers.

A "VIRTUAL" LABORATORY IN PROCESS CONTROL

In the field of industrial processes, use of the Internet for communication made its debut only recently [4]. The authors of the first such industrial communication tools pointedly stressed that with proper application, managers working in their offices in one country can at all times use the web browser to monitor the operational data from an electrolytic cell, a potroom or a smelter located in another country, and take the actions required, all this at the cost of a local telephone call. This amounts to a remote management information system, gathering and processing data for decision making purpose at managerial level.

The capacity of the web to host, communicate and share not only data, but also more sophisticated tools makes it possible to go further and to use the web for process simulation, model calibration, parameter adjustment and experimentation, analysis and control, not at the managerial level but at the level of the process itself.

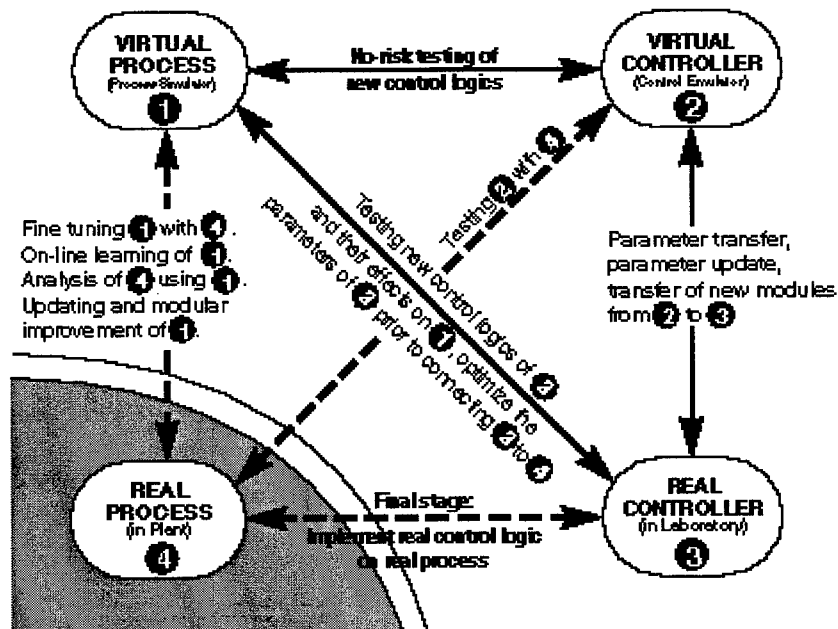


Fig. 2. The basic concept of a "virtual" laboratory for process control.

In recent years exploratory work was conducted at the Université du Québec à Chicoutimi in close collaboration with the aluminum industries, mainly Alcan International Limited's Arvida Research and

Development Center (ARDC). The results obtained encouraged creation of a new Process Control Laboratory. The concept behind this facility is to network through the Internet, real and virtual components of process control systems. Figure 2 shows the basic arrangement. The laboratory is built on four main components namely the real and virtual process, the real and virtual controller. Three out of four of these components are located in the laboratory proper whereas the fourth, the real process, stays in the plant where it belongs. Networking between the plant and the laboratory is done through the Internet, with the relevant security measures to ensure controlled access, confidentiality, data communication safety, through an appropriate implementation of user identification, encryption-decryption, error detection and error correction procedures.

Figure 2 shows a number of possible applications through networking of two or more of the four basic components of the laboratory. Through coupling of the virtual process (the process simulator) with the real process (in the plant), one can calibrate the former by applying data of the latter. The former can exercise on-line learning based on the data of the latter. The virtual process can be used as tool for carrying out a process analysis on the real process. In the case of complex processes, simulators are often built in modular form, and individual modules can be updated or improved through the coupling of real and virtual process.

By way of coupling the virtual process (the process simulator) with the virtual controller (the control emulator), it is possible to carry out inexpensive and risk-free tests for new control strategies. By way of coupling the virtual process with the real controller, we can test new schemes implemented in the real controller, identify their impacts on the virtual process and optimize the control parameters of the real controller prior to coupling it to the real process. Through a coupling of the virtual controller with the real one, we can perform the transfer of parameters, the update of parameters as well as the transfer of the new modules that are developed experimentally on the virtual controller for the purpose of incorporation into the real controller. Clearly, the ultimate aim is to couple the real controller with the real process, and the importance of this final stage justifies the many preparatory steps.

A PROGRESS REPORT

In recent years we have built a knowledge base for supervision in off-line consulting mode, of the feeding of the aluminum electrolytic cell [2]. Models were also built to simulate dynamic behavior of various types of cells [5]. The models are based on mathematical representation of energy and material balances, complemented by relations describing the chemical reactions, the physical properties and the operational parameters. These models serve as tools for research, process analysis and personnel training. Control emulators were also built to emulate the basic control actions generated by the cell's own control system.

Neural networks were built for the purpose of predictive and adaptive control of the cell [6, 7]. Predictive control aims at predicting the behavior of the cell and adjusting the feeding actions in anticipation to avoid or minimize the anode effects, whereas in adaptive control, the neural network recognizes the cell's characteristic curve and therefrom deducts the present thermal state of the cell, to which the control actions must be adapted.

Beside the electrolytic cell, our work covers a wide range of primary aluminum industry processes. As examples, a control system for controlling a rotary calcining kiln for petroleum coke, based in part on fuzzy logic, was patented in 1996 [1]. A dynamic model incorporating a control emulator was built in 1998 for the simulation of the calcining furnace for anthracite [8].

Setting-up of the new Process Control Laboratory is presently underway. Its conception is guided by the Virtual Laboratory concept presented earlier, and once operational, its short term mandate will be to experiment and implement control strategies based on computational intelligence, some of which have been developed or investigated earlier. Beside the process models used as "virtual processes", the laboratory will be equipped with hardware and software tools, comprised mainly of PID tuning analysis and simulation software, PLC CPU's and modules, operator interface and expert system tools. As a first stage, two complete control systems are built and they have been chosen to be representative of a major part of process controls in the primary aluminum industry. The first is an electrolytic cell control system and the

second is control of a liquid metal furnace, which could be any of the many categories of furnaces used in the industry such as melters, holders, mixers, casting furnaces or recycling furnaces. The cell control system focuses on cell feeding and on the many facets of cell operation such as anode adjustment, anode change, metal tapping or the suppression of anode effects. The furnace control system aims at burner control, stack control, combustion chamber temperatures (gas, refractories) and metal temperature control.

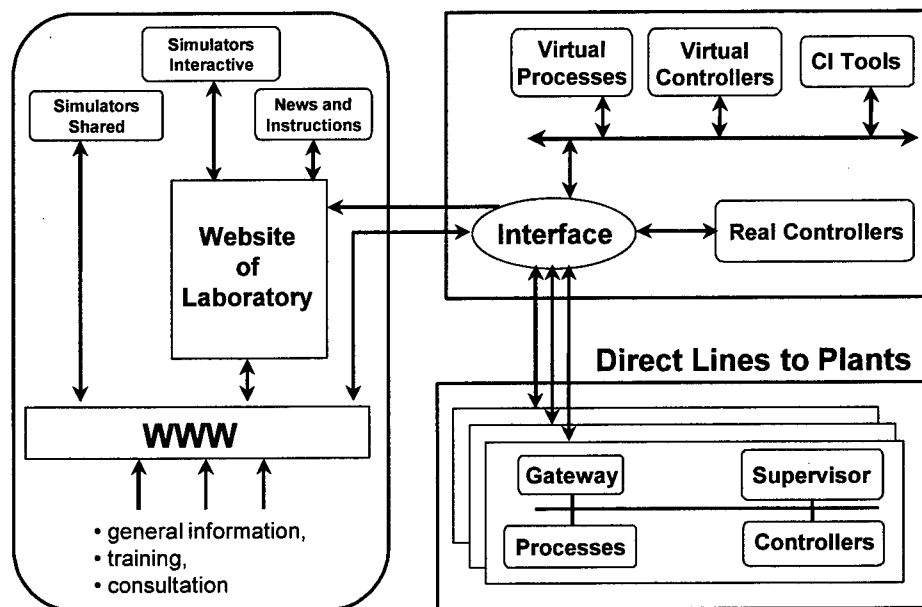


Fig. 3. The general configuration and networking of the process control laboratory.

Each of the two control systems is implemented according to the general configuration of Figure 3. In the laboratory, the virtual process, the virtual controller and the real controller together with the CI tools are connected together and to the "outside world" through an interface similar to the ones used in industry. The gateway to the plants is equipped with appropriate security measures to ensure controlled access and confidential data exchange. The laboratory components including the virtual process (the simulator) can be made available on the web in interactive mode (the webpage browser can interact with the simulator), or in shared mode, in which the visitor can use the simulator to carry out experiments and work together in real time with others at the laboratory.

CONCLUSION

Control of processes differs from control of the more homogeneous mechanical or electrical systems fundamentally by the fact that in the former, there is usually a higher number and variety of physical mechanisms involved, and not all of them can be formally understood or represented analytically in closed form. This is why direct human intervention through knowledge, experience, reasoning or induction plays a prominent role in process control. We are witnessing an ever-increasing capacity of computers to not only simulate complex physical processes but also duplicate human brain workings. On the other hand, new computer technologies are available to facilitate high-speed information and communication. This clearly provides unprecedented opportunities for a big leap forward in industrial process control with positive consequences on product quality and productivity.

ACKNOWLEDGMENTS

The research work on industrial process modeling and characterization was supported by the Industrial Research Chair in Process Engineering of the Université du Québec à Chicoutimi (UQAC), funded by the Natural Sciences and Engineering Research Council of Canada (NSERC), Alcan International Limited, and La Fondation de l'UQAC. The setting-up of the new Process Control Laboratory is funded by a major grant

to UQAC provided jointly by the Canadian Foundation for Innovation (CFI), the Ministry of Education of Quebec (MEQ) and Quebec's Aluminium Research and Development Center (CQRDA).

REFERENCES

1. J. Perron and M. Auger, 1996. "Process for Controlling Rotary Calcining Kilns, and Control System Therefor.", United States Patent Number 5523957.
2. C. Desbiens, 1992. "Base de connaissances pour la supervision de procédés." M. Eng. Thesis, Université du Québec à Chicoutimi.
3. "Aluminum Technology Roadmap Workshop", The Aluminum Association and The United States Department of Energy, Washington D.C. 1996.
4. M. De Souza Baltar, E.F.V. Da Silva, C.S. Filho, 1998. "Analyzing and Monitoring Aluminum Smelter Performance via Web Browsers.", J.O.M, 50(8), 14-16.
5. L. Tikasz, R.T. Bui, V. Potocnik, 1994. "Aluminum Electrolytic Cells : a Computer Simulator for Training and Supervision.", Engineering with Computers, 10, 12-21.
6. A. Meghlaoui, R.T. Bui, J. Thibault, L. Tikasz, R. Santerre, 1997. "Intelligent Control of the Feeding of Electrolytic Cells Using Neural Networks." Metall. and Materials Trans, 28B, 215-221.
7. A. Meghlaoui, R.T. Bui, J. Thibault, L. Tikasz, R. Santerre, 1998. "Predictive Control of Aluminum Electrolytic Cells Using Neural Networks." Metall. and Materials Trans, 29B, 1007-1019.
8. R.T. Bui, R. Hachette, G. Simard, J. Perron, J.F. Dessureault, 1999. "Computer Simulation of the Anthracite Calcining Furnace.", 1999 TMS Annual Meeting, San Diego, CA, February.

Modelling of the Flow Stress using BP Network

Y. Y. Yang, D. A. Linkens

Department of Automatic Control and Systems Engineering
University of Sheffield
Mappin St., Sheffield S1 3JD, UK

ABSTRACT

This paper addresses the development of a back-propagation neural network model for flow stress prediction based on plane strain compression test data. Basic concepts of the neural network modelling are given, followed by discussions on training data requirements and other critical issues in neural network modelling. Original training data have been obtained via many PSC tests for a low carbon steel (C430). Data pre-processing is very important in neural network modelling, especially when the data are from industrial processes where various disturbances are very likely. A two-stage data pre-processing procedure was proposed to deal with the PSC data: data rationalising and data filtering. The quality of the training data is significantly improved after the data pre-processing. The developed BP neural network model had been implemented on a Pentium-based personal computer. Simulation results show that the average output prediction error by BP network is less than 4% of the prediction range. The training error gradually decreases with increasing hidden neurons. However, increasing hidden neurons do impose a danger of over-training, with the validation error increasing instead of decreasing. Compromising between the training error and validation error, we suggest that a BP neural network with a single hidden layer and 10-20 hidden neurons should be sufficient for flow stress modelling.

INTRODUCTION

The accuracy of numerical simulation and many other design calculations (such as the rolling force, etc.) depends on the description of mechanical properties of the deformed materials. The strain hardening functions relating the yield stress to the temperature, strain and strain rate are commonly used in finite element models [1]. How to obtain an accurate strain-stress relationship becomes critical to the correct calculation of the finite element model. When the temperature and strain rate can be varied, strain-hardening function is not easy to obtain, whether from a physical-based or from a recursive model due to the high nonlinearity and complicated interaction among the stress, strain, temperature and strain rate. Many empirical functions have been proposed to calculate the stress, usually considering the effects of strain, temperature and strain rate [2, 3]. Recently, artificial neural networks have become a popular tool for flow stress modelling [4, 5], including the internal variable model considering the dislocation density as an extra input [6]. But the accurate modelling of the strain harden function for varying temperature and strain rate has not yet been fully achieved.

In this paper, we will focus on the modelling of the flow stress (determined by the strain hardening) using a back-propagation (BP) neural network. The training data are obtained through a series of laboratory plain stress compression (PSC) tests conducted on low carbon steel C430. A brief description of the PSC tests and the actual data collection is given, followed by data pre-processing treatment which improves the quality of the original data before they are actually fed to training the neural network. The neural network modelling procedures are then developed and tested, and typical results are given. Finally, concluding remarks, along with a brief discussion on further research work are outlined.

PSC TEST AND TRAINING DATA GATHERING

In order to investigate the flow stress behaviour under hot plane strain compression, a series of compression tests with different temperature, strain, and strain rate ranges, have been conducted on a hot PSC apparatus as shown in Fig. 1(a). Specimen (taken from low carbon steel C430) are prepared, with the geometrical

shape sketched in Fig. 1(b). Typical data for the specimen are $b_0 = 50$ mm, and $h = 10$ mm, with the b_1 , b_2 , b_3 and h affected by the total strain and the strain rate for the PSC test. The steel specimen is first heated in the heating furnace until it reaches the specified temperature, then it is put into the testing apparatus that conducts the pressing with the prescribed nominal strain rate. During the compression, specimen temperature (obtained via the embedded thermocouples), the displacement, and compression force are recorded. Strain, stress, and strain rate are later calculated from these recording data, with some necessary corrections for the geometry change of the specimen in three dimensions (width changing) [7].

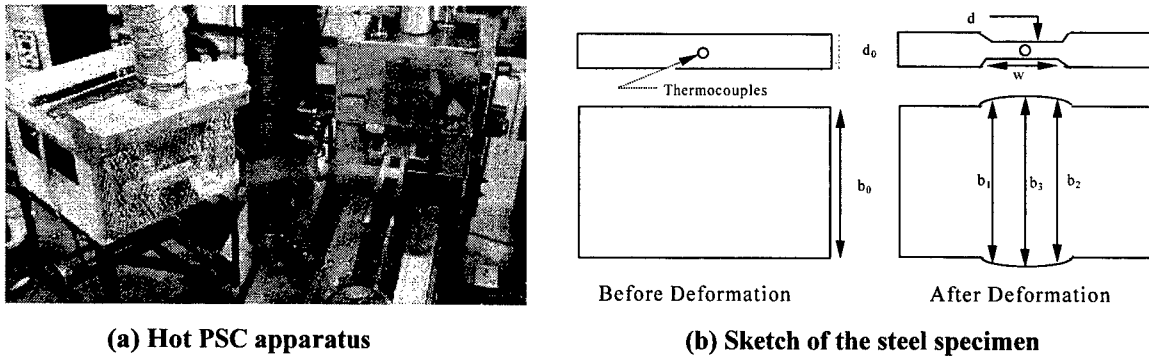


Fig. 1. Hot plain strain compression apparatus

39 SPC tests have been conducted for investigating the stress-strain relationship during the compression. The nominal PSC test condition range: nominal strain rate varies at 1, 10, and 50 s^{-1} , nominal specimen temperature varies at 900, 1000, and 1100°C , and maximum strain of about 2. Original recorded variables during the test are the displacement, the press load, and the specimen temperature. Other parameters, such as the initial width and thickness of the specimen, the final spread parameters (b_1 , b_2 , b_3) are measured either before or after the test. The strain, strain rate (which fluctuates around its nominal value), and stress are calculated from the above measurements, with compensations for the specimen dimension change and non-pure plain strain [7]. For the modelling work described here, strain, strain rate and specimen temperature, are chosen as inputs, with the flow stress as output.

DATA PRE-PROCESSING FOR THE PSC DATA

Data pre-processing is very important in any kind of black box modelling, especially when the original data are from real industrial processes where various disturbances are likely to intrude. After preliminary analysis of the PSC data, it revealed that the temperature measurements for some PSC tests are quite noisy and need proper treatment before they can be useful for modelling. The noise within the temperature measurement may be caused by electrical/magnetic disturbances around the PSC test apparatus. The strain rate also seems a little erratic, which might be caused by the derivative calculation of the strain within the calculation package. Moreover, the strain-stress data recorded at the beginning and end zones of the compression test are much more problematic than those in the middle of the test, due to the inaccurate zero calibration and the dramatic changes of strain rate and stress around the beginning and end zones.

After visualisation of the PSC test data (using Microsoft Excel), we designed a two-stage data pre-processing procedure. The first stage is referred to as *data rationalising*, in which data points at the beginning and finishing zones are cut out, and data points which are reckoned as irrational are deleted. The criteria of determining whether a data point is rational or not are mainly from the understanding of the PSC test process. For the specimen temperature, for example, judgement can be made by checking whether the temperature is in the feasible range and whether the temperature change between any two adjacent measurements is more than a possible value (threshold). The threshold values to determine whether a temperature is too low or too high, and whether a temperature change between two adjacent measurements are beyond reality come partly from experience of the PSC test, and partly from common sense. Since there is a second stage data pre-processing to follow, the threshold values are not so critical if they drift from

their idealised values. In our case, the temperature range is set to 700C (minimum) and 1250C (maximum), and the temperature changing threshold is set to 80C.

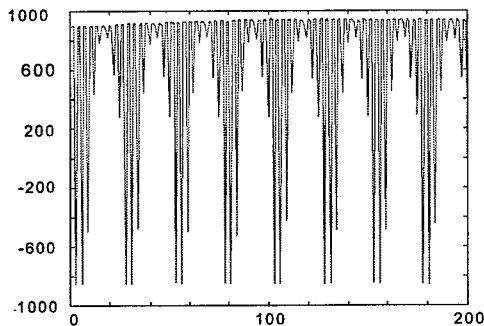
The second stage of data pre-processing is referred to as *data filtering* which tries to smooth out the noise still contained in the PSC test data after the first stage of the data pre-processing. Two kinds of filtering algorithms i.e., mean average filtering given by equation (1) and median average filtering given by equation (2), have been used here. We found that the median average filtering is better and more robust for the PSC test data encountered here.

$$y_f(k) = \frac{1}{2h+1} \sum_{i=k-h}^{k+h} y(k) \quad 1.$$

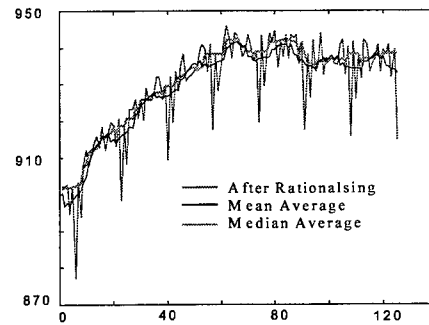
$$y_f(k) = \text{median}(y(k-h), y(k-h+1), \dots, y(k+h)) \quad 2.$$

where h is a non-negative integer which represents the half-zone (the half length of the data series to be averaged) of the mean or median filter, $y(k)$ is the original value of the variable at sample k , and $y_f(k)$ is the filtered value of the variable at sample k . The significant of the filter is controlled by the half-zone parameter h . General speaking, for a mean average filter, a large h will increase the smoothness of the filtered data, at a cost of low sensitivity. However, this does not hold for the median average filtering. If $h = 0$, no filtering is carried out.

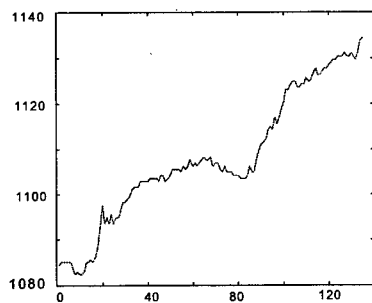
Fig. 2 shows two examples of the data pre-processing for the specimen temperature. The first data contains significant noise in its temperature measurement, as shown in Fig 2 (a-b), while the second contains less noise (as shown in Fig 2 (c-d)). The half-zone parameter is set to $h = 4$.



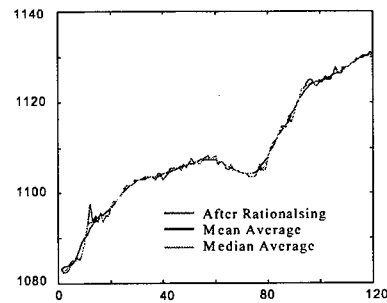
(a) Original measurements (Dat430si)



(b) After data pre-processing



(c) Original measurements (Dat430sq)



(d) After data pre-processing

Fig. 2. Data pre-processing for specimen temperature measurements

From the above examples we find that both the mean and median filtering are very efficient at removing the high frequency noise, and drive the filtered signal towards its true process values. However, the median filtering is much more robust than the mean filtering, since the outlier will have little contribution in the median filtering while in the mean filtering the outlier is weighted on the final filtered signal. If the original signals contain little noise, then the filtered signals (from either mean average filter or median average filter) will be very close to the original signals, as can be seen from Fig. 2 (d).

BP NETWORK MODELLING

In this paper, we use the BP neural network to model the flow stress behaviour during the PSC test. The BP network is the most commonly used neural network related to modelling, consisting of an input layer, several (typical 1 or 2) hidden layers, and an output layer [8]. Fig. 3(a) shows the structure of a three-layer BP network. The power of the BP network has been demonstrated by a number of workers, and research has indicated that a BP neural network has the potential to approximate any continuous non-linear function with arbitrary accuracy provided that there are enough hidden neurons [9].

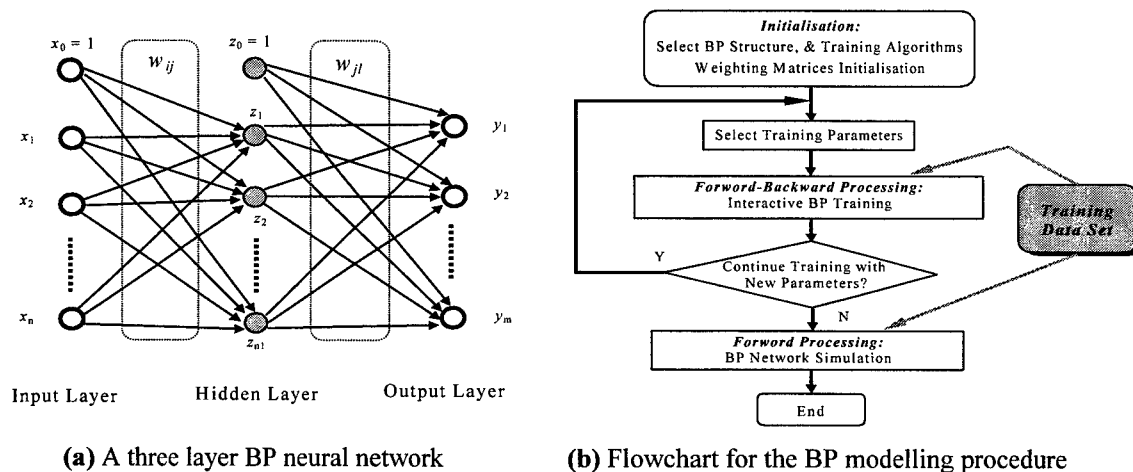
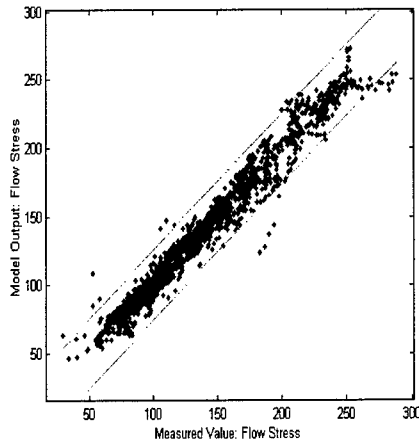


Fig. 3. A three layer BP network

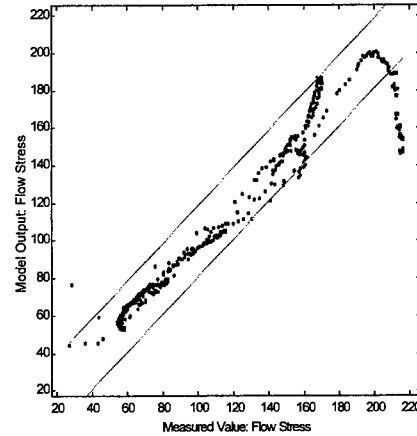
In BP neural network modelling, three stages are commonly involved, i.e., the initialisation, the forward processing, and the backward processing. The initialisation stage sets up the neural network architecture (layers and number of hidden neurons in each hidden layer), determines the activation function for the hidden layer(s), selects the training algorithms and parameters, and initialises the weighting matrices and bias with 'small' random values. The forward processing is to calculate the network outputs when presented with the input using the current network parameters. The backward processing is responsible for the training of the network (by adjusting its weight matrices) based on the error index (a measure representing the distance between the network output y and the desired target output y_d). One common algorithm for backward processing is the back-propagation of the error through the network to determine the updated weighting matrices and bias. In this paper, the Levenberg-Marquardt optimisation algorithm was used for the training the neural network [10]. The flowchart for the BP modelling procedure is shown in Fig. 3(b). The BP network modelling procedure is implemented under the Matlab environment. All the available strain-stress data are first combined into an overall data set, except 5 PSC tests which are deliberately kept out to form a testing data set. The overall data set is then divided into a training data set and a validation data set, with the amount of the data to be included in the training set controlled by a training ratio parameter *RatioT* (1-100%). The pattern of selecting the data from the overall data set to form the training and validation sets is controlled by a partition parameter *Mixed*, with options available for separate, sequential, or random partition.

Numerical simulations with various training algorithms and training parameters have been carried out on a Pentium 150 PC computer. Extensive simulations have been carried out to study the BP modelling capabilities, including the influence of hidden neurons, the effect of data partition and data filtering, etc.

Some typical simulation results are shown in Fig. (4-5). For the simulation presented here, the activation function for the hidden layer is sigmoid function *logsig*, while a pure linear function is used for the output layer. Input variables have been filtered by median filtering before being fed to the neural network, with the half-zone parameter $h = 4$. Data partitions use *Ratio = 50* and *Mixed* for sequential partition.



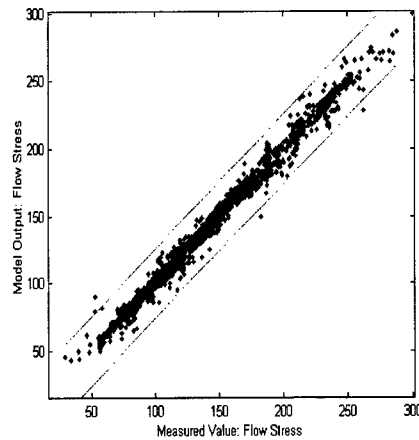
(a) Model prediction for the training data



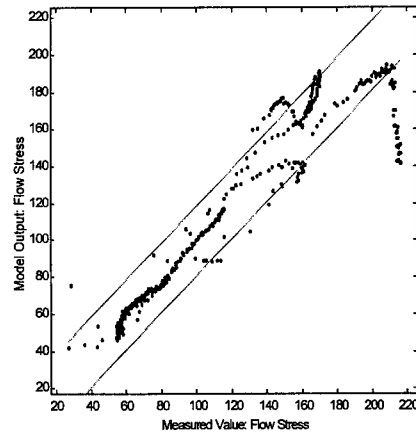
(b) Model prediction for the testing data

Fig. 4. BP network with 10 hidden neurons.

Training Error: 0.025711, Validation Error: 0.025749, Testing error: 0.036624



(a) Model prediction for the training data



(b) Model prediction for the testing data

Fig. 5. BP network with 30 hidden neurons.

Training Error: 0.014044, Validation Error: 0.014340, Testing error: 0.043212

DISCUSSION AND CONCLUSIONS

From the above results, we see that most of the output (flow stress) predicted by the BP network are within 10% of the error band (indicated by the straight lines in Fig. 4 and 5). The errors in Fig. 4 and 5 are normalised mean absolute error and the performance for training validation data are very close in all results, implying that the information contained in the test data is sufficient to predict the validation data. This confirms that sequential partition of the PSC test data into training and validation data sets with RatioT = 50 is satisfactory. The training error gradually decreases when hidden neurons are increased, although reduction is not dramatic over the range of 5 to 30. However, increased hidden neurons do impose a danger of over-fitting, with the validation error increasing instead of decreasing. For example, with 30 and 10 hidden neurons, although the training error dropped to 0.0140 from 0.0257, the validation error increased to 0.0432 from 0.0362 respectively. If we look at the testing data, the BP network with 30 hidden neurons (Fig 5 (b)) has more prediction points falling outside the 10% boundary than does its 10 hidden neurons counterpart (Fig 4(b)). This is a clear indication of over-training. Compromising between training error and validation error, the simulations suggest that 10 to 20 hidden neurons are sufficient for PSC flow stress modelling.

For testing data never seen by the BP neural network, the network gives a reasonably accurate prediction. However, the data points at the far right of Fig 4(b) are problematic even with the best trained BP network. When we checked these data, we found that most are from the same PSC test (Dat430cf) where temperature variation during the test was much higher than other tests. It can be argued that in Dat430cf, there may be incorrect measurements and hence, it should be dropped from the testing data set. If treated this way, the testing error is compatible with the training error and all the testing data fall into the 10% error band.

This initial work shows that BP neural networks can model complex flow stress behaviour provided sufficient training data are available. Next, we intend to develop a more general network able to model the stress-strain relationships of a wide variety of steels (not just a single steel). It is critical to find representative training data with rich-enough information for the class of steels in question. Since the data set will be large, data reduction and feature extraction will be required. After establishing a generalised BP neural network, it can replace the empirical flow stress model within the finite element model framework

ACKNOWLEDGEMENT

The assistance of Dr. Bruce Davenport in providing the PSC test data is gratefully acknowledged. Financial support from EPSRC is also acknowledged under Grant GR/II/73585

REFERENCES

1. O. Wiklund, 1996. "Rolling force models for temper rolling using finite elements and artificial neural networks", Proc. Steel Strip 96, Sept., Opava, Czech Republic.
2. J. H. Beynon, and C. M. Sellar, 1992. "Modelling microstructure and its effects during multipass hot rolling", ISIJ International, 32(3), 359-367.
3. H. Shi, A. J. McLaren, C. M. Sellars, R. Shahani, and R. Bolingbroke, 1997. "Constitutive equations for high temperature flow stress of aluminium alloys" Materials Science and Technology, 13, 210-216
4. Y. J. Hwu, Y. T. Pan, and J. G. Lenard, 1996. "A comparative study of artificial neural networks for the prediction of constitutive behaviour of HSLA and carbon steels", Steel Research, 67(2), 59-66.
5. Y. J. Hwu, and J. G. Lenard, 1996. "Application of neural networks in the prediction of roll force in hot rolling", Proceedings of the 37th MWSP Conference (ISS), Vol. XXXIII, 549-554.
6. J. Kusiak, and M. Pietrzyk, 1997. "Artificial neural networks applied as a history dependent constitutive model for hot forming of steels", Australia Pacific Forum on Intelligent Processing and Manufacturing of Materials, Gold Coast, Australia, July 14-17.
7. N J Silk, and M R Winden, 1998. "Interpretation of hot plane strain compression testing of aluminium specimens", IMPPETUS Report 003, June.
8. P. S. Neelakanta, 1994. Neural network modelling : statistical mechanics and cybernetic perspectives, Boca Raton, Florida; London : CRC Press.
9. K. Hornik, M. Stinchcombe, and H. White, 1989. "Multistage feedforward networks are universal approximators", Neural Networks, 2, 359-366.

10. C. M. Bishop, 1995. *Neural Networks for Pattern Recognition*, Clarendon Press, Oxford.

Intelligence in Materials Science II

The Heredity and Control of Microstructures of Liquid Metals during Rapid Cooling Processes

Rang-su Liu^{*}, Ji-yong Li^{}, Hai-rong Liu^{**}**

^{*} Department of Physics, Hunan University, Changsha, 410012, China

^{**} Department of Chemistry, Hunan University, Changsha, 410012, China

ABSTRACT

An investigation on the heredity and control of the microstructures of liquid metals during rapid cooling processes has been performed under different conditions by molecular dynamics simulations. From the simulations some important results have been obtained. First of all, in the system of liquid metal, which atom becomes the central atom of a small cluster is accidental and random. However, as long as a small cluster has been formed with a given atom as the center and some surrounding atoms, the cluster could be repeated with the same central atom and the same surrounding atoms again and again during the runs. The clusters would possess relative stability during the isothermal processes and heredity (continuity) during the rapid cooling processes. The stability and heredity of icosahedral clusters can be expressed quantitatively by their lifetime or repeatable times. In general, the lifetimes of these clusters increase with decreasing temperature, especially, below the glass transition temperature T_g . The number of clusters having a longer lifetime also increases with decreasing temperature. These will give us a new way to understand and control the heredity and transition mechanisms of microstructures of liquid and solid metals.

INTRODUCTION

Recently, many workers devoted themselves to the research of improving the macroscopic properties of metals and alloys. However, the macroscopic properties are mainly determined by their microstructures, and the microstructures are mainly determined by cooling processes from liquid metals and alloys. In order to improve their macroscopic properties, it is necessary for us to understand the relationships between the microstructures of liquid state and that of solid state for metals and alloys, especially the transition features of the microstructural configurations during their cooling processes. As we known, it is difficult to complete a tracking study for the transition processes of microstructures of liquid metals. With rapid development of computer technique, we can make such a simulation study for the transition processes of microstructures by means of molecular dynamics method. In recent years, some important results have been obtained in the authors' previous works^[1-5]. Especially, the heredity of microstructures of liquid metals is very interest for materials scientists since it may be play an important role during the transition processes of microstructures of liquid metals. Therefore, a deep research is worth making for understanding the physical origins of heredity and controlling its concrete process.

Based on the authors' previous works^[1-5], the main purpose of this paper is to study in detail the heredity of the microstructures of liquid metal Al by tracking its rapid cooling processes under different initial states using molecular dynamics method. From the simulation results, a clear picture was obtained to show that how the atoms in liquid metals gather to form some clusters and how the clusters further evolve to form some new type of clusters during their cooling processes. And the ways, for how to control the transition direction, especial the hereditary direction of some clusters, were discussed in detail from different purposes.

SIMULATION CONDITIONS AND METHODS

In this paper, as shown in Ref. [1-5], a molecular dynamics simulation study on the microstructure transitions of liquid metal Al during the rapid cooling processes has been performed under different initial states. All the simulations are made with the same system consisting of 500 Al atoms placed in a cubic box and run with periodic boundary conditions. The interacting interatomic potential adopted here is the

effective pair potential function of the generalized energy independent nonlocal model-pseudopotential theory developed by S. Wang et al.^[6-7], and the function is

$$V(r) = (Z_{\text{eff}}/r) [1 - 2/\pi] \int_0^\infty dq F(q) \sin(qr)/q \quad 1.$$

where Z_{eff} and $F(q)$ are, respectively, the effective ionic valence and the normalized energy wave number characteristic, which have been previously defined in detail. [6,7]. The pair potential is cut off at 20.0 a.u. (atomic unit) as shown in Fig.1. The time step is 10^{-15} s.

The simulations are started at $T=943\text{K}$ (which is 10K higher than the melting temperature T_m of metal Al). First, the system is run at this temperature, respectively, for 2500, 5000, 7500, 10000 and 15000 time steps to obtain five different initial states. For each initial state, let the temperature of the system decrease with the cooling rate of $33.5 \times 10^{12} \text{ K/s}$ from 943K to 50K. The atomic configurations are recorded at some particular temperatures during the cooling processes. Another run of 4000 time steps at each corresponding temperature is performed to obtain 20 different configurations for five corresponding initial system. Then detect the bond-type indexes between the related atoms using the Honeycutt-Andersen(HA) bond-type index method^[8]. Finally, analyze and compare the changes of the relative number of bond-types in the system to obtained some new results.

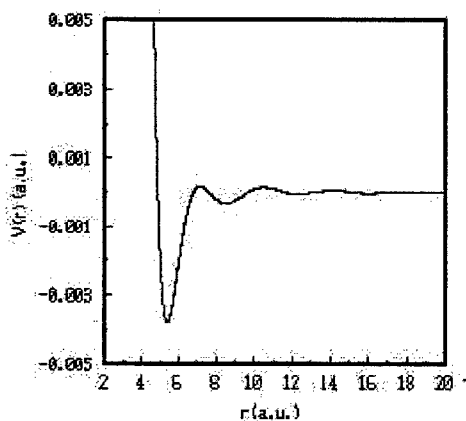


Fig.1. Effective pair potential $V(r)$ of liquid metal Al at 943K

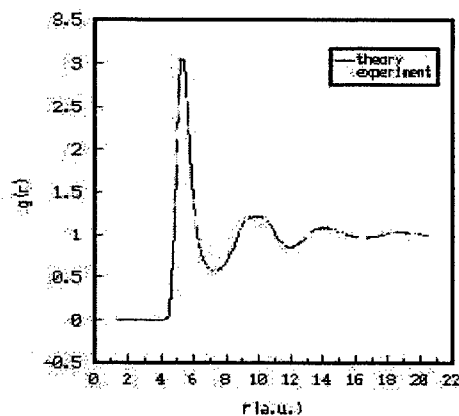


Fig.2. Pair relative distribution function of liquid metal Al at 943K.

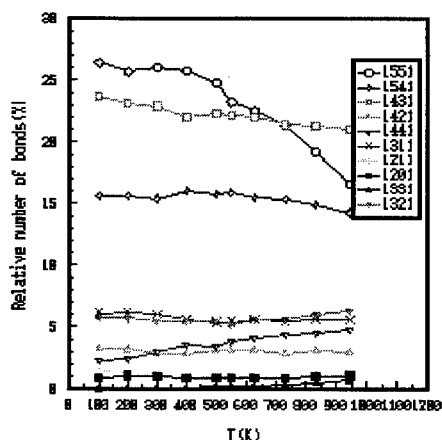


Fig.3. Relations of the relative number of various HA bond-types with temperature during rapid cooling process of liquid metal Al.

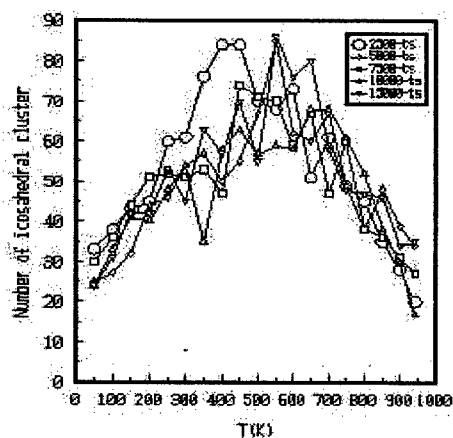


Fig.4. Relations of the number of icosahedral clusters for five different initial states with temperature during rapid cooling

SIMULATION RESULTS

At first, in order to detect the confidence level of the simulation results, we compare the pair distribution function $g(r)$ of the system obtained from the simulations with the experimental results given by Waseda^[9] and find that they are consistent very well each other, as shown in Fig.2.

Secondly, from Fig.3, it can be seen that of all the Honycutt-Andersen (HA) bond-types, the 1551 bond-type, related to icosahedral cluster, plays a critical role during the microstructure transition processes of liquid metal Al with decreasing temperature. For convenience of discussion, we choose the 1551 bond-type as the representation of all the bond-types so as to explain the main characteristics of simulation results.

From the recorded data of configurations deduced from initial states of 2500, 5000, 7500, 10000 and 15000 time steps, it can be clearly seen that at each given temperature, some icosahedral clusters can appear repeatedly with the same central atom and the same surrounding atoms as shown in Table 1. For simplicity, we use the repeatable times to express the lifetime of the clusters quantitatively during the isothermal and cooling processes. Go further, we also use the repeatable times to express the stability and heredity of icosahedral clusters quantitatively. With the decrease of temperature, the repeatable times of the icosahedral clusters are increased remarkably. For short, we only give the number of the central atoms that repeated more than 3 times and only the two results corresponding to the initial states of 2500 and 15000 time steps, as shown in Table 2 and Table 3.

Table 1.

The numbered central and surrounding atoms of icosahedral clusters repeated more than 10 times at 400K

No.of central atom		No.of surrounding atoms												
3	57	84	109	160	220	289	302	354	408	410	417	491		
149	10	116	132	160	192	320	326	347	419	455	464	497		
153	17	85	146	155	158	188	276	312	346	359	393	404		
212	34	71	100	135	167	179	208	259	267	291	361	416		
394	33	44	108	114	172	284	313	328	342	371	375	493		
431	18	78	87	120	138	189	194	204	215	296	363	365		
436	7	67	70	84	102	111	295	376	397	410	417	498		

DISCUSSION

From the simulations mentioned above, some important results can be seen as follows:

1. When the system under consideration is in liquid state or in solid state, during its isothermal runs, it can be clearly seen that which atom could become the central atom of a cluster and which cluster could be repeated, are entirely accidental and random. However, once an atom has become the central atom of a cluster, the cluster would appear with the same central atom and the same surrounding atoms again and again. The cluster, in fact, would possess relative stability and keep its initial configuration in all the isothermal runs at different temperatures and cannot be broken arbitrarily. Thus we can choose the central atom numbered as the label of a cluster in the system, for example, some repeatable icosahedral clusters such as 3, 149, 153, 212, 394, 431 and 436, are shown in Table 1. Those clusters are the results of initial state of 2500 time steps and they can appear repeatedly up to 20 times as shown in Table 2.

As we know, during all the runs, it can be seen that once a cluster has formed, that would either keep its central atom and the same surrounding atoms for all the runs, or be broken. However, in any case, it cannot be seen that the cluster is formed again with the same central atom and different surrounding atoms. We think this is a significant property of the microstructures of metals. This stability gives us an important basis to understand the physical origins and the concrete mechanism of the microstructure transitions of liquid metals during rapid cooling processes.

Table 2. Relationship of the various numbers of icosahedral clusters with temperature during rapid cooling process of liquid metal Al (with initial equilibrium time of 2500 steps)

Temperature (K)	No. of icosahedral structures	No. of repeatable icosahedral structures	No. of unrepeatable icosahedral structures	Repeatable times with central atoms numbered in brackets
943	20	6	14	
900	28	3	25	
850	35	6	29	3(243);
800	46	18	28	3(57,205,376,400,449,498);
750	49	19	30	4(145,348);3(20,172,189,262,396);
700	61	25	36	9(262); 8(278,447); 6(198); 5(8,294); 4(266); 3(19,34,46,129,199,245,387,461);
650	51	25	26	7(237); 6(152,354,416); 5(34,55,220); 4(328,401,418); 3(265,317,324,327,404,410,499);
600	73	30	43	9(184); 7(151,448); 6(83,171,285,416); 5(68,154,395); 4(41,199,238); 3(13,73,91,211,270,337);
550	68	38	38	14(420); 11(162,384); 10(360); 8(55,270,439); 7(487); 6(494); 5(65,73,81,370,480); 4(58,75,259,371); 3(189,215,261,331,378,419);
500	70	48	22	9(24,44,73); 8(250); 7(42,79); 6(1,356,382); 5(96,122,129,244,249,262); 4(32,127,158,196,301,343,360); 3(34,48,52,121,270,312,332,349,358,371,428,433);
450	84	53	31	12(234,314,399,477); 11(73,413); 10(357); 9(28,206,356); 8(1,63,162,461); 7(131,328); 6(482); 5(129,209); 4(11,29,164,166,244,266,269,276,481); 3(7,8,32,86,127,130,153,158,315,420,424);
400	84	53	31	17(461); 16(314); 14(131); 12(86,455,497); 11(422); 10(159,481); 9(203,267); 8(123); 7(261,322); 6(49,60,318,376); 5(190,409); 4(55,160,171,264,361,466); 3(74,234,259,270,321,393);
350	76	47	29	18(73,409); 17(261); 16(376,422,466); 15(171,284); 14(81); 13(75,267); 10(369); 9(7,33,456); 7(383,392); 6(82); 5(437); 4(94,225,226,321,396,398,412); 3(96,98,219,228,271,314,362,400,463,473);
300	61	45	16	18(284); 16(210,76); 15(463); 14(431,477); 13(219); 11(75,271); 10(73,394,400); 9(78); 8(314,437,456); 7(261,422,494); 6(43,240,408,486); 5(72,315,339,398,409); 4(329,351,466); 3(30,203,225,293,300);
250	60	44	16	20(284,394); 19(431); 15(422); 13(78,314); 12(219,486); 11(43,300); 10(153,240,437,477); 9(261,376,400); 8(30,149); 7(132,179,271,362,484); 6(25,73,210,339,351); 5(361,463); 4(3,225,494); 3(41,112,171,408);
200	45	38	7	20(339,394,431,437); 19(149,300,351,436,484); 18(477); 14(84,486); 13(153,219); 12(3,284); 11(73,462); 10(210,433); 9(271); 7(212); 6(78,314); 5(240,342); 4(75,132,261); 3(43,361,408,435,493);
150	43	40	3	20(149,153,212,361,394,436,437,486); 19(431); 18(351,462,477); 16(339); 15(300); 13(325); 12(449); 11(84); 10(3,73,210); 7(219,250,342,362); 6(311,314,408); 5(132,148,261,370); 4(41,43,435); 3(78,240,293,433);
100	38	31	7	20(3,149,153,394,431,436,437,462); 19(477,486); 18(212,325,351); 15(339); 13(300); 11(84,132); 10(408); 9(210,342); 8(261,362,449); 5(219,311); 4(240,361); 3(271,314,433);
50	33	31	2	20(3,149,153,212,394,431,436,437,462,477,486); 19(325,351); 17(408); 16(449); 15(339); 14(84,210,300); 10(342,362); 9(261); 6(219,311,314,361); 5(240); 3(132,433);

Table 3. Relation of the various numbers of icosahedral clusters with temperature during rapid cooling process of liquid metal Al (with initial equilibrium time of 15000 steps)

Temperature (K)	No. of icosahedral structures	No. of repeatable icosahedral structure	No. of unrepeatable icosahedral structures	Repeatable times with central atoms numbered in brackets
943	34	5	29	3(418);
900	39	6	33	4(66);
850	48	7	41	5(277); 3(148);
800	38	12	26	5(418); 4(310); 3(420);
750	49	16	33	4(300,451); 3(121,168,220,252,475,491);
700	67	34	43	7(392); 5(61,285); 4(88,145,322); 3(181,182,239,298,333,364,415);
650	60	31	29	7(454); 6(44); 5(97,211,301,409); 4(200,227,255,379,460,483); 3(40,221,275,386);
600	62	31	31	9(108); 8(28,360); 6(14,172,363,398); 5(116,178); 4(23,234,342,466); 3(45,289,336,440);
550	85	38	43	12(437); 11(118,301); 10(69); 9(287); 8(196,442); 6(8,25,223,476); 5(90,248,386,400,408); 4(78,84,208,273,296,349,379); 3(109,151,163,213,238);
500	67	39	28	11(435); 10(118); 8(25,109,271); 7(145,185); 6(43,355,448,496); 5(235,424); 4(24,318,375,411,461,463); 3(7,75,186,304,313,406,471,479);
450	55	38	17	19(108); 18(316); 13(192,463); 10(26,355,399); 9(186); 8(231); 7(67,145,264,335,437); 6(74,341,487,490); 5(43,235,272,440); 4(257,338,476); 3(16,27,144,194,237,282,360,499);
400	50	36	14	18(355); 16(399); 14(304); 13(339); 12(342); 11(31); 10(144,316); 9(13,237); 6(108,239,264,437); 5(21,118,257,467); 4(192,282,360); 3(135,150,178,228,328);
350	57	38	19	19(355); 18(21); 17(261,381); 16(43); 15(237,342); 14(144,228,335); 13(331); 11(36,316,437); 10(95,377); 8(99); 7(26); 6(69,186); 5(17); 4(271,298,400); 3(13,195,200,284,291,374,399);
300	54	42	12	20(21); 19(43,355,400); 17(13,496); 16(17,284); 14(36,335); 12(26,261); 11(65,95,144,331); 10(342); 9(237,437); 8(254,316,458); 7(314); 6(151); 5(129,204,377,447); 4(271,282,445); 3(186,192,228);
250	46	32	14	20(17,21,284,355); 19(314,316,447); 18(26); 17(36,101); 13(320,400); 10(204); 9(43,271); 8(399,493); 7(13,95); 6(347); 4(151,254); 3(108,336,367,437);
200	44	29	15	20(17,21,26,36,284,314,316,355); 19(447); 17(320); 15(342); 13(74); 12(43,101,400,493); 11(271); 10(95,399); 9(291); 6(108,239); 5(347); 4(91); 3(144,370);
150	32	27	5	20(17,21,26,36,271,284,291,314,316,355,399,447); 18(320); 12(95); 11(342); 10(73); 9(400); 7(108,239,493); 6(101); 5(467); 4(386); 3(43,498);
100	27	23	4	20(17,21,26,36,284,291,314,316,320,355,399,447); 17(271,342); 16(95); 13(400); 10(239); 9(101); 7(73); 6(43); 5(493); 4(108); 3(386);
50	25	22	3	20(17,21,26,36,271,284,291,314,316,320,355,399,447); 17(95); 16(73,239); 15(342); 13(400); 11(386); 8(43); 7(493);

2. During rapid cooling, with a decrease in temperature, we can see that: the total number of the icosahedral clusters in all the systems deduced from different initial states are increased rapidly. For instance, the numbers of clusters deduced from an initial state of 2500 time-steps, as shown in Table 2 and Fig.4, are increased at first from 20 to 84 in the range of 943K – 450K, then through a maximum of 84 smoothly in

the interval of 450 - 350K corresponding to the glass transition temperature T_g , thereafter, decreased rapidly from 84 to 33 in the range of 400- 50K. For the numbers of clusters deduced from initial state of 15000 and other time steps, there are the similar ways as shown in Table 3 and Fig.4. But the size and positions of their maximum for each system are different, namely, the corresponding glass transition temperature T_g are different and it can be moved with different initial states.

The total numbers of repeatable icosahedral clusters in all the systems deduced from different initial states are also almost rapidly increased at first, then through the maximum 53 and 42 of themselves, respectively, for the initial states of 2500 and 15000 time steps, and decreased rapidly as shown in Table 2, 3 and Fig.4.

The repeated times of the repeatable icosahedral clusters in all the systems as above-mentioned are also increased, although the highest repeated times of the clusters not increasing gradually, sometimes they are up or down abruptly. As the temperature downs to below T_g , especially to 250~300K, the highest repeated times will be and keep the saturation value of 20 times and can not be changed again. However, the number of the icosahedral clusters repeated 20 times will be greatly increased from 1 to 11 and 13, respectively, for initial states of 2500 and 15000 time steps.

It is interesting that the total numbers of the non-repeatable icosahedral clusters in all the systems as above-mentioned are also increased as shown in Table 2 and 3. As temperature downs, respectively, to 600K and 550K, the maximum of each system is the same value 43. In this case, it can be seen that almost 50~60% of the total number of icosahedral clusters are the non-repeatable. And then it is decreased rapidly to 2 ~ 3 at 50K, that is to say, only 5 ~10% of the total number of icosahedral clusters can not be repeated, and the 90 ~95% of them can appear repeatedly, namely, most of them are very stable. These results are just the expected from thermodynamics theory.

3. It can be clear seen that the another important result is the heredity of the icosahedral clusters during the rapid cooling processes. In general, it is difficult that a cluster can be repeated in isothermal runs at some given temperature and thereafter can be repeated again in the next isothermal runs at another temperature. However, a few clusters can keep their repeated, namely, they possess continuity or heredity, in the next several isothermal runs. For example, the clusters labeled 477, 314 and 271 appeared and repeated more than 6 ~ 10 times in the range of 500K ~50K, for almost 8 ~ 10 temperature intervals, as shown in Table 2. Similarly, the clusters labeled 271, 284, 355 and 447 also appeared and repeated more than 7 ~ 10 times in the range of 500K ~50K, for almost 8 ~ 10 temperature intervals, as shown in Table 3. Therefore, we can consider that these clusters (477, 447, 355, 314, 284 and 271) possess higher stability and heredity (continuity) during the cooling processes.

Especially, it is worth notice that the clusters labeled 271, 284 and 314 can appear and to be repeated in the two systems deduced from initial states of 2500 and 15000 time steps. Some clusters, such as the cluster 437, can appear in the range of 350K ~ 50K in Table 2, but it only appears in the range of 550K ~ 250K, and can not appear in the lower temperature of 200K ~ 50K in Table 3. In addition, the cluster 449, it can appear in the range of 850K ~ 800K, disappear in the range of 750K ~ 200K, and then appears in the range of 150 ~ 50K in the system as shown in Table 2. However, just this cluster 449, it only appears at 800K, and can not be found again in the range of 750K ~ 50K in the system as shown in Table 3. These results tell us that under different initial conditions, the clusters possess different level of stability and heredity.

Go further, from the results mentioned above, it can be clearly seen that all the clusters possessing stability and heredity can be divided into three levels: higher, middle and lower. The clusters 271, 284, 314 and 437 can be considered as the higher level, since they can appear in the two systems more than 10 ~ 15 times. The clusters 149, 153, 210, 394 and 431 only appear in the system more than 5 times as shown in Table 2. The clusters 21, 26, 36, 316 399 and 447 also only appear in the system more than 5 times as shown in the Table 3., therefore, those can be considered as the middle level. And the most clusters listed in Table 2 and 3 can appear 1 ~5 times during all the cooling processes, those can be considered as the lower level.

4. From these results, it is demonstrated that the stability and heredity of the clusters are different for different initial conditions and can be controlled by changing their initial conditions (including temperature, pressure, cooling rate, etc.) This will give us a new way for the microstructure design of metallic materials

CONCLUSIONS

From the results and discussions mentioned above, we have obtained a very clear picture about how the metal atoms gather to form clusters and how the clusters further evolve to form some new types of clusters during the rapid cooling processes.

1. It can be clearly seen that which atom could become the central atom of a cluster and which cluster could be repeated, are entirely accidental and random. The clusters can appear again and again with the same central atom and the same nearest neighboring atoms, but cannot appear with the same central atom and not the same surrounding atoms, until they are dissociated completely.
2. The stability and heredity of icosahedral clusters can be expressed quantitatively by their lifetime or repeatable times. Above the glass transition temperature T_g , the stability and heredity of clusters in the system under different conditions are almost in the similar level.
3. However, below the glass transition temperature T_g , the stability and heredity of clusters in the system have remarkable variations, some clusters could appear in all the cooling processes. With the decrease of temperature, the number of the clusters having more heredity is also increased remarkably.
4. The stability and heredity of the clusters are different for different initial conditions and can be controlled by changing their initial conditions (including temperature (as shown in this paper), pressure, cooling rate, etc. (will be shown in other papers)). This will give us a new way to design the microstructure of metal and alloy materials.

ACKNOWLEDGMENT

This work was supported by the National Natural Science Foundation of China. The authors gratefully acknowledge the support of K. C. Wong Education Foundation, Hong Kong.

REFERENCES

1. R.S. Liu, D.W. Qi and S. Wang, Phys. Rev. B, 45, 1992, 451- 453.
2. R.S. Liu and S. Wang, Phys. Rev.B, 46, 1992, 12001- 12003.
3. R.S. Liu, J.Y. Li, and Q.Y. Zhou, Chinese Science Bulletin, 17, 1955, 1429- 1433
4. J. Y. Li, Z. Zhou, R. S. Liu, Q. Xie and P. Peng, J. Mater. Sci. Technol. , 14, 1998, 461-464.
5. R. S. Liu, J. Y. Li, Z. Zhou, K. J. Dong, P. Peng and Q. Xie, Mater. Sci. Eng., B57, 1999, 214-217
6. S.Wang and S.K.Lai, J. Phys. F, 10, 1980, 2717- 2737.
7. D.H.Li, X.R.Li and S.Wang, J.Phys.F, 18, 1986, 309 -321.
8. J.D.Holnecutt and H.C.Andersen, J. Phys. Chem, 91, 1987, 4950 - 4963.
9. Y.Wasweda, The Structure of Non-Crystalline Materials, McGraw, HJill, New York, 1980, p20.

Artificial Intelligence Approach to the Internal Variable-based Rheological Model For Steels

J. Kusiak, M. Pietrzyk

Akademia Gorniczo-Hutnicza, Mickiewicza 30, 30-059 Krakow, Poland

ABSTRACT

The paper is a continuation of the authors' earlier work dealing with application of artificial neural networks to the prediction of yield stress in hot forming of metals. At present, the task of the network is to predict a time-derivative of the dislocation density during hot deformation. The inputs are the state of the material defined by the current dislocation density and by the time-integral of strain, the current strain rate and temperature. The flow stress curve is determined from the dislocation density vs. strain function, which is calculated using a finite difference technique in which the time-derivative of the dislocation density is supplied by the artificial neural network. Examples of calculations are presented for the axi-symmetrical compression of low carbon steel.

INTRODUCTION

The accuracy of numerical simulation of metal forming processes depends strongly on the description of mechanical properties of the deformed material. The predictive capability of the model describing the flow stress during hot deformation can be improved when the state of the material is related to "so-called" internal variables. In hot forming, these variables comprise the dislocation density, the recrystallised volume fraction and the grain size. Earlier research aimed at developing an internal variable model for microstructure evolution in steels [1,2]. The main difficulty in practical application of this model is related to evaluation of material constants. It is shown in [2] that an application of the inverse technique to evaluation of these constants often presents serious problems. The cost function is flat with local minima and searching for a global minimum is time consuming. Moreover, it is often difficult to find the constants, which give proper results of simulation in a wide range of temperatures and strain rates.

Thus, an attempt was made to apply artificial neural network to predict the influence of various components of the internal variable model on the overall behaviour of the material [3]. The basic assumption of the approach is that the time-derivative of the dislocation density is the only output parameter of the neural network. The state of the material, represented by the current dislocation density, is the input parameter of the network. This approach gave good results when single strain rate tests were investigated [3]. The main objective of the current work was to extend the analysis to different strain rates. Training of the network was done using experimental data from results of axi-symmetrical compression tests performed at three temperatures and three strain rates. The trained network was implemented into a finite-element code and simulation of the tests was performed.

CONVENTIONAL ANN YIELD STRESS MODEL

Artificial neural networks have become a powerful tool in simulation and control of various processes. Numerous examples of an application of the ANN in metal forming can be found in the scientific literature. Among the many publications, those dealing with control of rolling mills [4,5] as well as with prediction of yield strength in plate mills [6], rolling loads [7,8,9], plate bending in asymmetrical rolling [10] and roll bending in 4-high stands [11] should be mentioned. Prediction of a material's resistance to deformation is one of the fields in which the ANN technique is very useful. In the conventional application of ANN to the modelling of the yield stress, the inputs are temperature, strain rate and strain while the output is yield stress. Typical results obtained by the authors from artificial neural network are presented in [11]. Good agreement between measured and predicted yield stress was obtained and it was concluded that artificial neural network

is able to reproduce stress-strain curves with a peak and plateau, which are characteristic of dynamic recrystallisation. On the other hand, the neural network, which was trained using current temperature, strain rate and strain as inputs, maintained all the drawbacks of the conventional stress-strain equations obtained by an approximation of experimental data, like for example, the Voce equation [12]. In the following section, the suggestion of a new approach is presented to apply ANN to predict yield stress.

INTERNAL VARIABLE ANN MODEL

The main assumption in the internal state variable model is that the evolution of stress during plastic deformation is governed by the evolution of dislocation populations [1]. This leads to a concept that hardening is controlled by a competition of storage and annihilation of dislocations, which superimpose in an additive manner. Since the mechanical strength of the obstacles to dislocations is related to the dislocation density, the yield stress accounting for a softening is calculated as:

$$\sigma = \mu b \rho^{0.5} \quad 1.$$

where: b = Burgers vector; μ = shear modulus; ρ = dislocation density

The dislocation density in this model should not be treated as an average value. Rather, the entire spectrum of the dislocation densities must be considered, as shown in [1,2]. Evaluation of the material constants in this model often presents serious difficulties. Application of the ANN to predict the dislocation density during hot plastic deformation allows us to avoid these difficulties [3]. In this approach, the current time-derivative of the dislocation density Φ is an output parameter. The inputs are the state of the material described by the current average dislocation density (ρ) and by the time-integral of the strain rate (φ), the current strain rate ($\dot{\epsilon}$) and temperature (T). As a consequence, when the ANN supplies the time-derivative of the dislocation density, the flow stress curve is calculated from equation (1) using a finite difference technique to determine the average dislocation density:

$$\rho_{i+1} = \rho_i + \Phi \Delta t \quad 2.$$

where: t - time, Φ - time derivative of the dislocation density calculated by the ANN as a function of current dislocation density, strain rate, temperature and time integral of strain rate:

$$\Phi = \frac{d\rho}{dt} = F(\rho, \varphi, T, \dot{\epsilon}) \quad 3.$$

Contrary to conventional approaches, the present model of the flow stress is an incremental type. The artificial neural network is used in each time step of the simulation.

RESULTS

Experimental Procedure

The tested material was a carbon-manganese steel containing 0.22%C, 1.26%Si, 0.016%P, 0.03%S, 0.1%Cr, 0.09%Ni, 0.27%Cu and 0.003%Al. All tests were performed on the deformation dilatometer DIL 805. The axi-symmetrical samples measuring 5 mm in diameter and 10 mm in height were preheated at 1150°C for 10 min, cooled in the furnace to the test temperature and compressed. The tests were performed with a constant die velocity and an average strain rate was calculated for each test. Current temperatures and compression loads were monitored.

Training the Network

Training was done using experimental data, which comprise the results of measurements of the compression loads at various test conditions. The experimental stress-strain curves were calculated from the measured loads accounting for the current contact area and for the influence of friction. These curves were differentiated

graphically and the time-derivatives of the dislocation density were calculated by reversing Equation 1. These derivatives were used for training the artificial neural network. Typical results obtained for the three test temperatures (1150°C, 1050°C and 950°C) and for the three strain rates (0.53 s⁻¹, 1.22 s⁻¹ and 2.8 s⁻¹) are presented in Figure 1. The solid lines represent experimental curves. It can be seen that agreement between the measurements and predictions is very good. It is also seen that the artificial neural network predicts correctly the stresses for temperatures of 1100°C and 1000°C, which are between those used in the testwork.

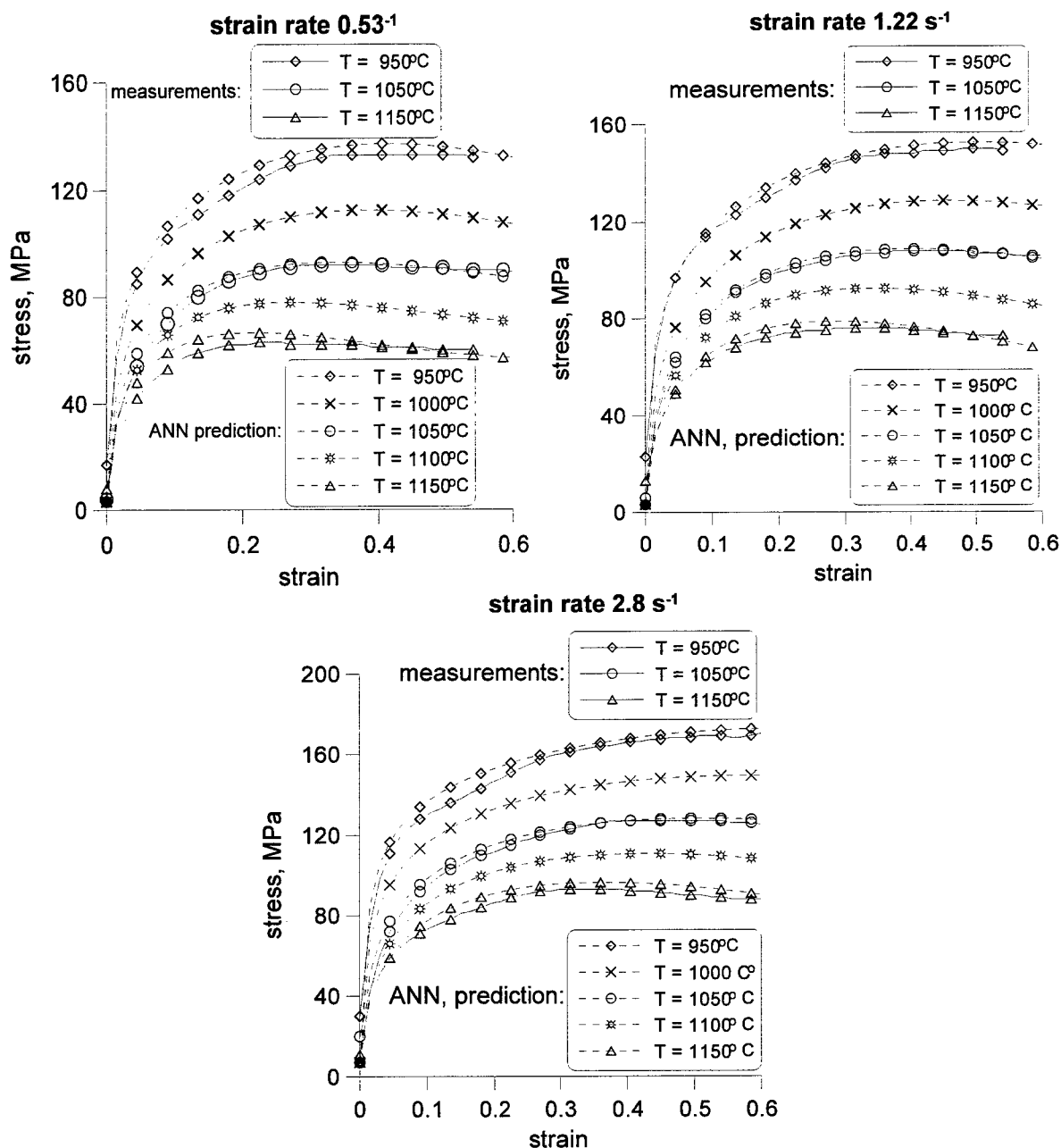


Fig. 1. Comparison of measured and calculated (ANN) stress-strain curves for different temperatures and strain rates.

The ability to predict flow stress variations during hot deformation at varying conditions is a potential advantage of the developed model, ensuing from its incremental character. Performance of the model for the deformation under varying strain rate is discussed below.

VARYING STRAIN RATE TESTS

The predictive capabilities of the model are tested by simulating varying strain rate processes. Typical results are presented in Figure 2. It was assumed that the strain rate changes rapidly between 0.53 s^{-1} and 2.88 s^{-1} . The change appears at a strain of 0.15, which is below the peak strain. The curves in Figure 2 represent calculations for constant strain rates of 0.53 s^{-1} and 2.88 s^{-1} , measurements for the constant strain rates of 0.53 s^{-1} and 2.88 s^{-1} and predictions for the strain rate changing rapidly between 0.53 s^{-1} and 2.88 s^{-1} . It is seen that after changing the strain rate, the predicted behaviour of the material does not reach the value determined by the new equation of state. This contrasts with the experimental observations for transient behaviour of C-Mn steels. Figure 3 shows typical measurements of stress-strain curves under varying strain rates obtained in [13]. The transient phase observed in the experiment and predicted by the internal variable model of [1,2] is shorter than that predicted by the current model. Conventional constitutive models do not predict transient behaviour at all (dotted line in Figure 3) so it can be concluded that the ANN trained using constant strain rate and temperature data fails in the situation of rapidly-changing conditions of forming. On the other hand, the ANN reproduces stress-strain curves for constant strain rates much better than does the internal variable model of [1]. Since variations in the strain rate in most metal forming processes are reasonably small, the rheological model based on the ANN can be efficient and useful in describing the yield stress in the finite element models.

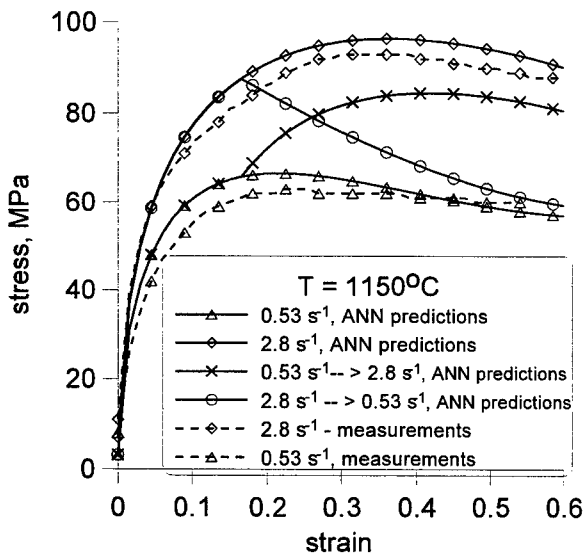


Fig. 2. Stress-strain curves calculated from the model for constant and varying strain rates compared with measurements for constant strain rates.

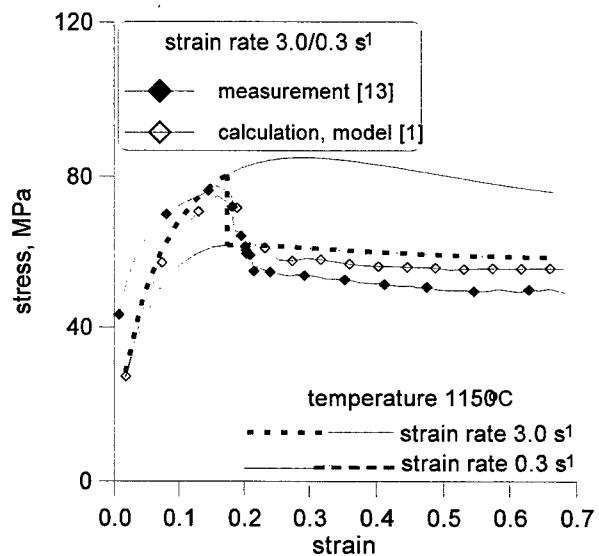


Fig. 3. Measured and calculated stress-strain curves for strain rate decreasing rapidly from 3 to 0.3 s^{-1} at a strain of 0.18 [13].

FEM SIMULATION

The developed ANN model is useful in the simulation of forming processes. Applying it as a constitutive law in the finite element approach has tested capabilities of the model. FEM program used in the calculations is described in [14]. It is based on the rigid-plastic flow formulation coupled with the solution of Fourier equation assuming Galerkin integration scheme. The ANN is used for determination of the yield stress in the Levy-Mises flow rule. Axisymmetrical compression was considered as an example. Figures 4 and 5 show typical results of calculations of strain field and yield stress field for the experiment carried out at temperature 1150°C and strain rate 0.53 s^{-1} . Calculated force as a function of height reduction is presented in Figure 6. The influence of increasing contact surface and the effect of friction can be well-understood in Figure 6. The force increases rapidly when the height of the sample decreases. All the results of FEM simulation agree with those obtained for the constitutive models based on closed-form equations obtained from an approximation of the experimental data (see [15]).

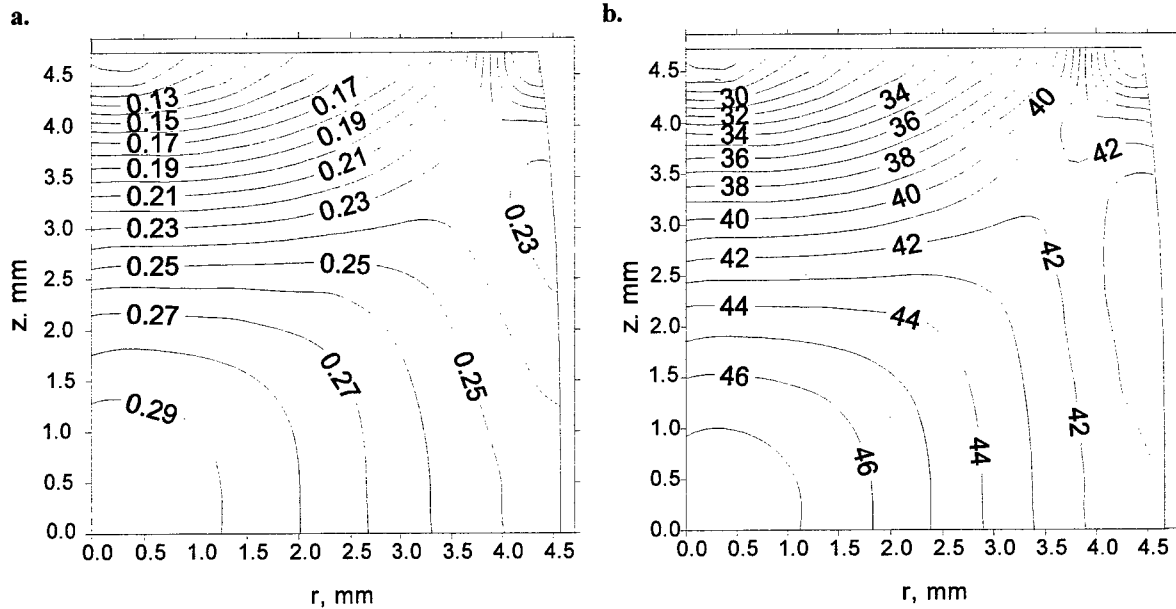


Fig. 4. Fields of: **a.** the effective strain and; **b.** the yield stress -- at the cross section of the sample deformed with a strain of 0.2; temperature of 1150°C, average strain rate of 0.53 s^{-1} .

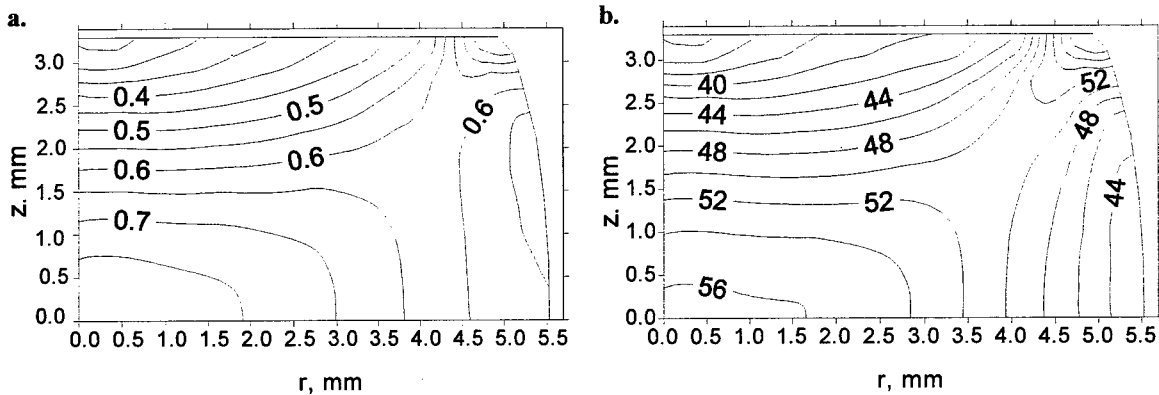


Fig. 5. Fields of: **a.** the effective strain and; **b.** the yield stress -- at the cross section of the sample deformed with a strain of 0.4; temperature of 1150°C, average strain rate of 0.53 s^{-1} .

CONCLUSION

An application of artificial neural networks to predict the yield stress in hot forming of metals has been presented in this paper. In this approach, the output of the ANN is a time-derivative of the dislocation density. The work is a continuation of the author's earlier research described in [3]. Influence of the strain rate has now been accounted for in the present work. The inputs of the network are the state of the material defined by the current dislocation density and by the time-integral of strain, the temperature and the current strain rate. Variations of average dislocation density as a function of strain are calculated using a finite difference technique. Good accuracy of the model was obtained for constant conditions of deformation.

The incremental technique is able to simulate processes involving varying strain rates and/or temperatures. However, as shown in Figure 2, the model fails to predict properly the transient behaviour of material when rapid changes in the strain rate appear. According to the predictions, after changing the strain rate, the stress never reaches the value determined by the equation of state for new conditions of deformation. This is in contradiction with experimental observations (see [13]). Thus, a conclusion can be drawn that the proposed model predicts perfectly the yield stress of steels during hot deformation only under reasonably stable

conditions. The results in Figures 4 and 5 show that the model can be implemented as a constitutive law into the FEM code. Application of this model to transient conditions requires further research.

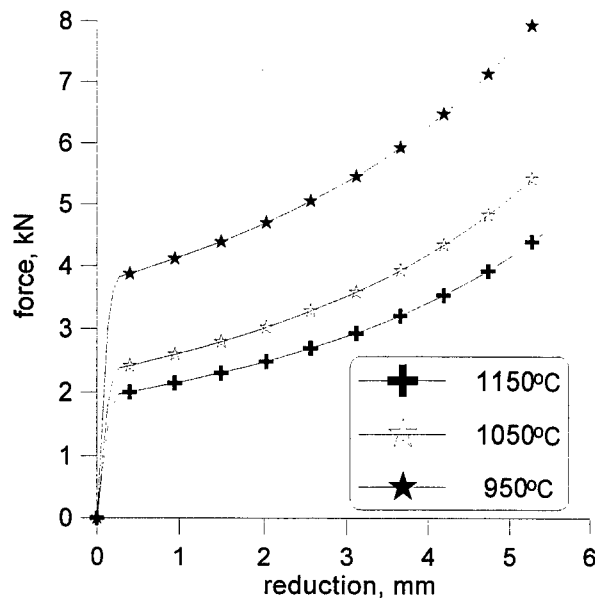


Fig. 6. Predicted force vs. reduction for various temperatures with an average strain rate of 0.58 s^{-1} .

ACKNOWLEDGEMENTS

Financial assistance of KBN (Grant No. T 08B 042 14) and NATO is gratefully acknowledged

REFERENCES

- Pietrzyk, M., 1994. Metall. Foundry Eng., 20, 429.
- Pietrzyk, M., Roucoules, C., Hodgson, P.D., 1995. Proc. NUMIFORM'95, (S.-F Shen, P. Dawson, eds) Ithaca, 315.
- Kusiak, J., Pietrzyk, M., 1997. Proc. IPMM'97, (T. Chandra, S.R. Leclair, J.A. Meech, B. Verma, M. Smith, B. Balachandran, eds), Gold Coast, 240.
- Roscheisen, M., Hofmann, R., Tresp, V., 1992. Advances in Neural Information Processing Systems 4, (M. Kaufman, ed.), 659.
- Too, J.J.M., Ide, K., Maheral, P., Pussegoda, N., Sherwood, E.G., Gomi, T., 1995. Proc. 37th Mechanical Working and Steel Processing Conference, Hamilton, 555.
- Tsoi, A.C., 1992. Advances in Neural Information Processing Systems 4, (M. Kaufman, ed.), 698.
- Y.J. Hwu, J.G. Lenard, 1995. Proc. 37th Mechanical Working and Steel Processing Conf., Hamilton, 549.
- Larkiola, J., Myllykoski, P., Nylander, J., Korhonen, A.S., 1996. Metal Forming96, (M. Pietrzyk, J. Kusiak, P. Hartley, I. Pillinger, eds), Krakow, J. Material Processing Technology, 60, 381.
- Wiklund, O., 1996. Steel Strip96, Opava, 136.
- Kusiak, J., Pietrzyk, M., Wilk, K., 1997. Proc. KomPlasTech97, (A. Piela, J. Kusiak and M. Pietrzyk, eds), Ustron-Jaszowiec, 207 (in Polish).
- Kusiak, J., Dudek, K., Svetlichnyj, D., Liszka, P. 1998. Proc. NEUROMET'98, (J. Kusiak, ed.), Krakow, 67 (in Polish).
- Hodgson, P.D., Collinson, D.C., 1990. Mathematical Modeling of Hot Rolling of Steel (S. Yue, ed.), Hamilton, 239.
- Pietrzyk, M., Kuziak, R., 1997. Proc. COMPLAS 5, (D.R.J. Owen, E. Onate, E. Hinton, eds), Barcelona, 1363.
- Pietrzyk, M., Lenard, J.G., 1991. Thermo-Mechanical Modelling of the Flat Rolling Processes, Springer-Verlag, Berlin.
- Majta, J., Lenard, J.G., Pietrzyk, M., 1996. ISIJ Int., 36, 1094.

The Mechanism of Electrolytic Al₂O₃ Coating on MAR-M247 Superalloy

S. K. Yen and C. C. Chang

Institute of Material Engineering, National Chung Hsing University,
250, Kuo-Kang Rd., Taichung, Taiwan 40242, R.O.C.

ABSTRACT

Through an analysis of cathodic polarization experiments in H₂O, NaNO₃, HCl and a mixture of HCl and NaNO₃, solutions respectively, the mechanism of coating electrolytic Al₂O₃ on MAR-M247 superalloy was investigated. We suggest that the cathodic polarization curves in Al(NO₃)₃ can be comminuted into 4 steps: 1. $H^+ + e^- \rightarrow H_2$ (-0.1V~-0.35V), 2. reduction of $Al^{3+}(H_2O)_3$ complex ion: $2Al^{3+}(H_2O)_3 \cdot xH_2O + 6e^- \rightarrow 2Al(OH)_3 \cdot xH_2O + 3H_2$ (-0.35V~-0.65V), 3. diffusion of $^{+}(H_2O)_3$ complex ion (-0.65V~-0.9V), and 4. reduction of H₂O: $2H_2O + 2e^- \rightarrow H_2 + 2OH^-$ (-0.9V ~ -4V). X-ray diffraction diagrams show the Al(OH)₃ gel transforms into amorphous-Al₂O₃ $\xrightarrow{623K}$ γ -Al₂O₃ $\xrightarrow{973K}$ δ -Al₂O₃ $\xrightarrow{1123K}$ θ -Al₂O₃.

Keywords: mechanism, electrolytic Al₂O₃ coating, superalloy

INTRODUCTION

Ceramic coatings seem ideal for use as high-temperature materials and in severely corrosive environments. Since superalloys are used in these extremes, ceramic coatings are often applied to enhance operation properties [1]. Forming ceramic coatings using an electrochemical method is a relatively new technique and has been used to deposit oxide coatings such as ZrO₂ on to metallic and non-oxide substrates [2-6]. Electrolytic deposition of aluminum hydroxide gel from an aqueous solution of aluminum nitrate (Al(NO₃)₃), followed by annealing, has been applied on SiC to increase substrate resistance to environmental attack such as high temperature oxidation [2]. In comparison with other deposition technologies such as chemical vapor deposition (CVD), physical vapor deposition (PVD) and plasma spraying, gel deposition has several potential advantages, including cheap deposition technology, energy economies, the ability to cover complex shapes and various materials, and application of multi-component oxides [7]. The sequence of reactions leading to alumina formation is considered to be the following [2]:

1. Dissociation of aluminum nitrate: $Al(NO_3)_3 \rightarrow Al^{3+} + 3NO_3^-$ 1.
2. Formation of hydroxide: $Al^{3+} + 3(OH)^- \rightarrow Al(OH)_3$ 2.
3. Dehydration to amorphous alumina: $Al(OH)_3 \rightarrow Al_2O_3 + 3H_2O$ 3.

Though the sequence of reactions leading to Al₂O₃ has been suggested, there is lack of analytical evidence to identify which cathodic reaction happens in the range of the applied cathodic voltage. Sometimes mud cracks and/or hydrogen-bubble effects in the coating, deteriorate film uniformity due to the uncertainty of either the electrochemical or drying mechanisms. Therefore exploring the precise electrochemical or drying mechanism should help us reach an optimal process to precisely control quality and quantity of the Al₂O₃ coating. In this study, we have carried out several cathodic polarization tests in H₂O (5.5 KΩ), HCl, NaNO₃, Al(NO₃)₃, and a mixture of HCl and NaNO₃ aqueous solutions respectively, to identify the regions to which these various reductive reactions belong.

EXPERIMENTAL

Sample Preparation

A MAR-M247 nickel base alloy was cut into discs with a diameter of 13 mm. All specimens were polished to a mirror finish with 1μm Al₂O₃ powder, then degreased by detergent and further ultrasonically cleaned in deionized water and acetone, then dried by N₂ gas. The nominal chemical composition of the MAR-M247 superalloy is given in Table 1.

Table 1. Nominal chemical composition of MAR-M247 superalloy.

Element	Cr.	Co	Mo	W	Ta	Al	Ti	C	B	Zr	Ni
Wt%	8.3	10.0	0.7	10.0	3.0	5.5	1.0	0.14	0.0015	0.05	Bal.

Polarization Tests

The MAR-M247 discs were electrochemically polarized in a naturally aerated 0.01 M $\text{Al}(\text{NO}_3)_3$ aqueous solution ($\text{pH} = 3.5$) by EG&G Princeton Applied Research 273A Potentiostat M352 software. Polarization tests were simulation in deionized water with an electrical resistance of 5.5 $\text{K}\Omega$ ($\text{pH} = 7.0$), HCl ($\text{pH} = 3.5$), 0.01 M NaNO_3 ($\text{pH} = 6.4$), and a mixture of HCl and 0.01 M NaNO_3 aqueous solution ($\text{pH} = 3.5$), respectively. The potential was swept from an initial potential of 0 V (AgCl) to a final potential of -4 V (AgCl), at a scanning rate of 1 mV/sec. The reactants of the simulated solutions are given in Table 2.

Table 2. The reactants of the simulated solutions.

Solution	Solution Reactants
H_2O	H_2O ; O_2 (8.2 ppm)
HCl	H_2O ; O_2 (7.8 ppm); H^+ ; Cl^-
NaNO_3	H_2O ; O_2 (8.2 ppm); Na^+ ; NO_3^-
HCl + NaNO_3	H_2O ; O_2 (8.2 ppm); H^+ ; Cl^- ; Na^+ ; NO_3^-
$\text{Al}(\text{NO}_3)_3$	H_2O ; O_2 (7.5 ppm); H^+ ; Al^{3+} ; NO_3^-

Electrolytic Deposition and Annealing

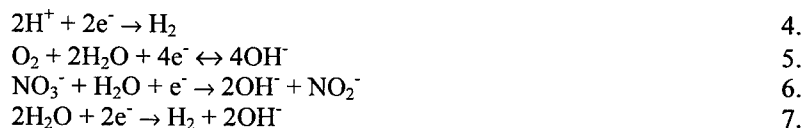
The electrolytic deposition of Al_2O_3 on MAR-M247 specimens was conducted in 0.01M $\text{Al}(\text{NO}_3)_3$ aqueous solutions at a voltage of -0.7V for 500s, using the same potentiostat. The alloy disc was the cathode, graphite was the anode and saturated AgCl served as the reference electrode. The above electrolytic conditions gave the most efficient deposition in our experiment. The specimens with $\text{Al}(\text{OH})_3$ gel coatings were then dried naturally in air and annealed in air at 473, 623, 973, and 1123 K for 2 hrs, respectively.

SEM and XRD

The surface morphology of coated and post-annealed specimens was observed by scanning electron microscopy (SEM, JEOL, JSM-5400, Japan). The crystal structure of the Al_2O_3 coating on the MAR-M247 substrate was analyzed by X-ray diffraction (XRD) in a MAC MO3X-HF Diffractometer, with Cu K α radiation ($\lambda = 1.5418 \text{ \AA}$), 2θ in the range $30^\circ - 85^\circ$, at a scanning rate of $4^\circ/\text{min}$, a voltage of 40 kV, and a current of 30 mA.

RESULTS and DISCUSSION

The cathodic polarization curve for 0.01 M $\text{Al}(\text{NO}_3)_3$ ($\text{pH} = 3.5$) is shown in Fig.1 -- curve $\text{Al}(\text{NO}_3)_3$. This curve can be divided into four steps. The first is -0.1 V ~ -0.35 V, the 2nd -0.35 V ~ -0.65 V, the 3rd -0.65 V ~ -0.9 V, and the 4th -0.9 V ~ -4 V. The cathodic reactions appropriate to the case of aqueous solutions of aluminum nitrate ($\text{Al}(\text{NO}_3)_3$) may be described as follows:



The reactants of 0.01M $\text{Al}(\text{NO}_3)_3$ should be H_2O ; O_2 ; H^+ ; Al^{3+} ; NO_3^- . Therefore polarization tests were simulated in pure water with an electrical resistance of 5.5 $\text{K}\Omega$ ($\text{pH} = 7.0$), HCl ($\text{pH} = 3.5$), 0.01 M NaNO_3 ($\text{pH} = 6.4$), and a mixture of HCl and 0.01 M NaNO_3 aqueous solution ($\text{pH} = 3.5$).

The cathodic polarization curves in pure H_2O and NaNO_3 as shown in Fig. 1, indicate the same limiting current density ($1.5 \times 10^{-5} \text{ A/cm}^2$). According to the reactants of H_2O and NaNO_3 , we suppose that the limiting current is due to the diffusion limit of O_2 in reaction 5.

The cathodic polarization curve in HCl has a prominent limiting current density of 5.5×10^{-5} A/cm². The reactants of HCl may be H₂O, O₂, H⁺, and Cl⁻. In comparison with H₂O and NaNO₃, the limiting current densities are different. The limiting current density of 5.5×10^{-5} A/cm² should be the diffusion limit of H⁺ in reaction 4., since it is much larger than the limiting current of O₂ in reaction 5.

The cathodic polarization curve in HCl + NaNO₃ shows a prominent limiting current density of 4×10^{-5} A/cm². This is very close to that of H⁺ above. Therefore, it is also considered as the diffusion limit of H⁺ in reaction 4. but with a slightly lower current density. This means that the addition of NaNO₃ reduces the diffusion limit of H⁺ in reaction 4., possibly due to the interaction between H⁺ and NO₃⁻.

When the applied voltage is more negative than -0.9 V, we noticed a lot of H₂ bubbles on the electrode surface for all solutions except pure water. This is due to the reduction of H₂O in reaction 7. The limiting current increased with decreasing electric resistance of aqueous solution because of the voltage drop of the solution between the two electrodes.

From the above analysis, the first step (-0.1 V ~ -0.35 V) with a limiting current density of 5×10^{-5} A/cm² in 0.01M Al(NO₃)₃ is considered to be the diffusion limit of H⁺ in reaction 4., since this current density is in the range of the diffusion limit current density of H⁺ [8]. However, this is not as obvious as that observed in ZrO(NO₃)₂ solution [9]. The fourth step (-0.9 ~ -4 V) is reduction of H₂O in reaction 1. However, the current density of the second step (-0.35 ~ 0.65V) in ZrO(NO₃)₂ solution is much more than the others. Compared with NaNO₃ solution, the only difference is Al³⁺ to Na⁺. Therefore, it is clear that the increased current density is due to Al³⁺ in the Al(NO₃)₃ aqueous solution. Possibly, formation of complex ion Al³⁺(H₂O)_{3+x} acts to reduce the activation energy of the reduction of H₂O via the following reaction:



Consequently, the third step (-0.65 ~ -0.9 V) with a current density of 1×10^{-3} A/cm², is the diffusion limit of Al³⁺(H₂O)₃. The process of electrolytic alumina coating on MAR-M247 superalloy as illustrated in Fig. 2 consists of: 1. dissociation of aluminum nitrate, 2. formation of complex ion Al³⁺(H₂O)₃, 3. diffusion and migration of Al³⁺(H₂O)₃, and 4. reduction of Al³⁺(H₂O)₃ and formation of Al(OH)₃ hydroxide.

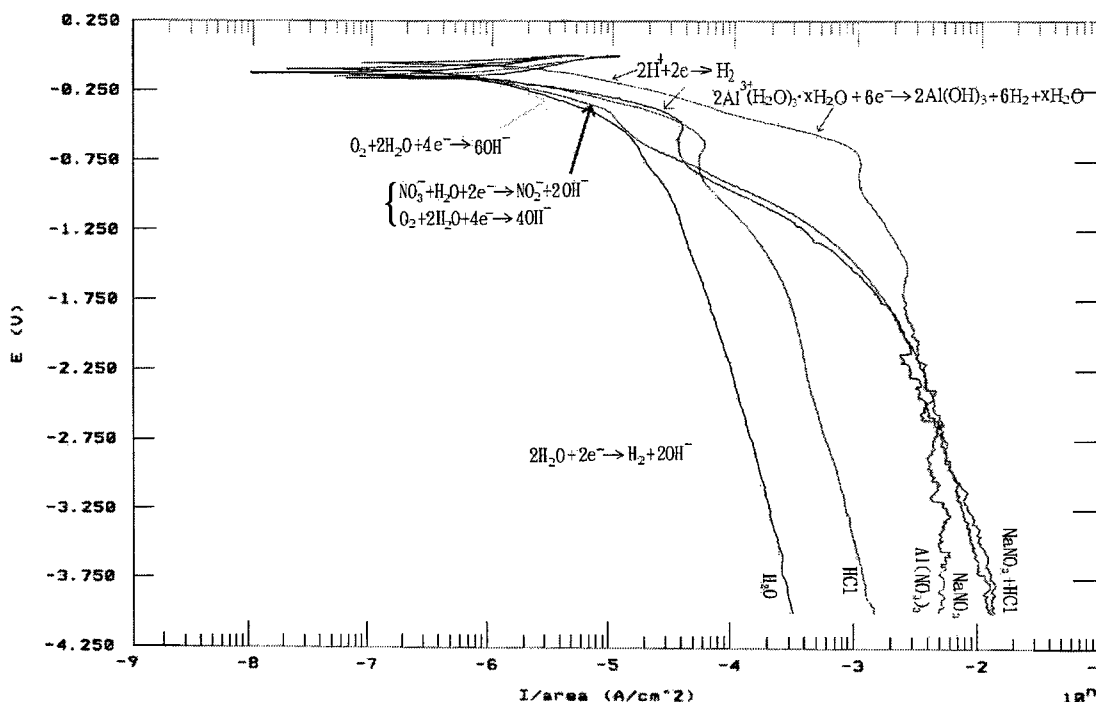


Fig. 1. The cathodic polarization curve in all of the reaction solutions.

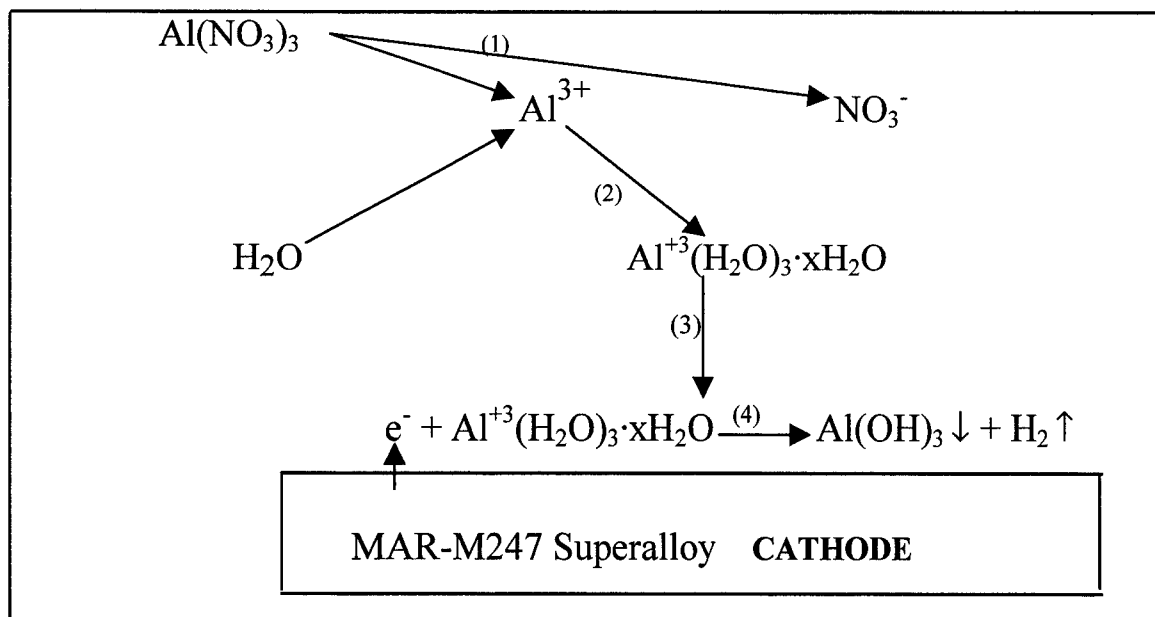


Fig. 2. The electrolytic deposition alumina diagram

The SEM micrographs of specimens coated at - 0.7 V for (a) 500 sec, and (b) 1000 sec after natural drying in air are shown in Fig. 3. Mud cracks increase with increased deposition time and the bubble effect can be found at a voltage of -1.0 V for 500sec (Fig. 3(c)), due to the reduction of H_2O in reaction 7. The best deposition voltage is during the third step (-0.65 ~ -0.9 V) which is considered to be reaction 8.

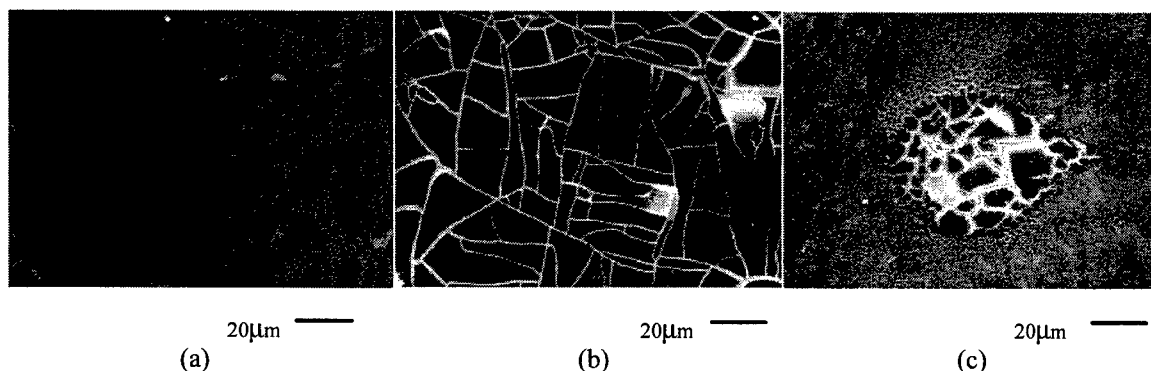


Fig. 3. SEM micrographs of specimens coated at - 0.7 V for (a) 500 s., (b) 1000 s. (c) at - 1.0 V for 500s. after natural drying in air (750 X)

The XRD patterns of specimens coated at - 0.7V for 500sec after natural drying in air, then annealing at 473, 623, 973 and 1123 K for 2hrs, for MAR-M247 superalloy are shown in Fig. 4. The diagrams show the $Al(OH)_3$ gel transforms into amorphous $\xrightarrow{623K} \gamma-Al_2O_3 \xrightarrow{973K} \delta-Al_2O_3 \xrightarrow{1123K} \theta-Al_2O_3$.

CONCLUSIONS

1. From the above discussion of the cathodic polarization tests, it is suggested that the cathodic polarization curves of MAR-M247 superalloy in $Al(NO_3)_3$ can be comminuted into four steps: 1. $H^+ + e^- \rightarrow H_2$ (-0.1V ~ -0.35V), 2. The reduction of $Al^{3+}(H_2O)_3$ complex ion: $2Al^{3+}(H_2O)_3 \cdot xH_2O + 6e^- \rightarrow 2Al(OH)_3 \cdot xH_2O + 3H_2$ (-0.35V ~ -0.65V), 3. The diffusion limit of $Al^{3+}(H_2O)_3$ complex ion (-0.65V ~ -0.9V) 4. The reduction of H_2O : $2H_2O + 2e^- \rightarrow H_2 + 2OH^-$ (-0.9V ~ -4V). The best deposition condition is at the third step.

2. X-ray diffraction shows $Al(OH)_3$ gel transforms into amorphous- $Al_2O_3 \xrightarrow{623K} \gamma-Al_2O_3 \xrightarrow{973K} \delta-Al_2O_3 \xrightarrow{1123K} \theta-Al_2O_3$.

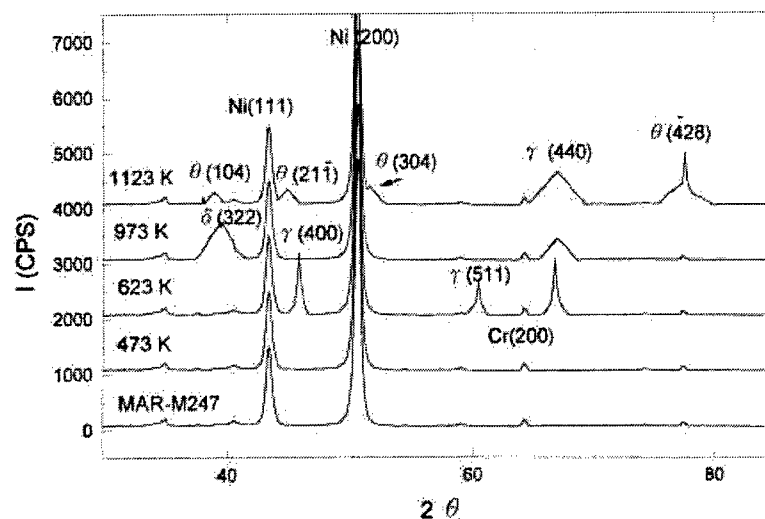


Fig. 4. XRD patterns of specimens coated after annealing and uncoated specimens.

REFERENCES

1. D. M. Comassar, 1991. Surface Coatings Technology for Turbine Engine Applications. *Metal Finishing*, 89(3), 39-44
2. R. Cham, G. Stark, L. Gal-Or and H. Bestgen, 1994. Electrochemical ZrO_2 and Al_2O_3 coatings on SiC substrates. *J. Materials. Science.*, 29, 6241-6248
3. L. Gal-Or, I. Silberman and R. Cham, 1991. Electrolytic ZrO_2 coatings (I). *J. Electrochem. Soc.*, 138, 1939-1942
4. R. Cham, I. Silberman and L. Gal-Or, 1991. Electrolytic ZrO_2 coatings (II). *J. Electrochem. Soc.*, 138, 1942-1946
5. R. Cham, I. Zhitomirsky, L. Gal-Or and H. Bestgen, 1997. Electrochemical Al_2O_3 - ZrO_2 coatings on non-oxide ceramic substrates. *J. Materials. Science.*, 32, 389-400
6. S. K. Yen and T. Y. Huang, 1998. The characterization of the electrolytic ZrO_2 coating on Ti-Al-4V. *Materials Chemistry and Physics*, 6, 214-221
7. R. G. Biswas, J. L. Woodhead and A. K. Bhattacharaya, 1997. Corrosion studies of inorganic sol-gel alumina coatings on 316 stainless steel. *J. Mater. Sci. Lett.*, 16, 1628-1633
8. Denny A. Jones, 1996. *Principles and Prevention of Corrosion*, 2nd Ed., 97
9. S. K. Yen, 1997. Modeling Electrolytic ZrO_2 Coating on Ti. 192nd Meeting of The Electrochemical Society, Paris, France, 97-2, 636

Automated Stress Control of Electroplated Nickel-Phosphorus Alloy

George Yu, Martin Williams, Tzu-Chern Horng, I. B. Huang

Department of Mechanical Materials Engineering
National Huwei Institute of Technology
Huwei, Taiwan 632

ABSTRACT

A process for automatic control of internal stress in electroplated nickel-phosphorus alloy using controlled current electrodeposition on a conductive substrate from a single electroplating bath yields a multiple-layered coating. This multi-layered deposit is a sequence of two alternating layers. One of these layers is a phosphorus-rich coating which is characterized by compressive stresses while the other layer is a tensile-stressed coating of low phosphorus content. The apparatus employed to achieve this automated control of internal stress includes a non-contact linear sensor for monitoring stresses in the coating and a programmable power supply for the current source.

INTRODUCTION

These internal stresses are often of little importance in deposits of weak and ductile metals such as tin, lead and cadmium. However, in stronger more brittle deposits such as nickel the internal stress is of great importance. In electroforming operations as well as subsequent mechanical operations stress may cause severe distortions on the workpiece. Excessive stress can also cause premature failure by accelerating corrosion or by decreasing fatigue strength. The magnitude of the intrinsic internal stress produced in electrodeposits is dependent upon bath compositions and operating conditions. Various organic additives are incorporated into plating baths to act as "stress relievers". They are chosen so as to produce low tensile-stressed deposits or compressively-stressed deposits. This is because compressive stress is normally less detrimental than tensile stress as it does not have a tendency to lift the coating from the substrate.

Current density is also an important factor in relation to the internal stress of an electrodeposited material. For example, the variation of composition of a nickel-phosphorous alloy is primarily a function of current density. It illustrates a decrease in the phosphorous content of a nickel-phosphorous alloy with increasing current density [1]. The internal stress of the alloy is in turn related to phosphorous content in the coating, as illustrated by Ref. [2,3,4,5,6]. At a phosphorous content of approximately 11% or higher the internal stress is neutral or compressive. Below this amount the internal stress in the deposit is tensile. So by using a controlled current and hence controlled current density it is possible to control the internal stress in the electrodeposition.

EXPERIMENTAL WORK

This study relates to a process for automatically controlling the stress in a coating using a controlled current mechanism. Therefore, This method eliminates the need for stress relief operations. It also eliminates the need for organic additives used as stress relievers. This reduces organic contamination in the plating bath and also eliminates the need for removal by attendants. Using a nickel-phosphorous alloy as a test case, the process for automatic stress control produces a multi-layer deposit. This multi-layer coating is a sequence of two alternating layers. One layer is a phosphorous rich coating, characterized by compressive internal stress, while the other layer is a tensile stressed layer having a low phosphorous composition.

Figure 1 shows a coating analysis apparatus for analyzing the internal stress of a coating. A testing cell (1), with a fixture (3) and a metal plate (4) mounted therein, is utilized to electroplate a specimen (6). One surface of the specimen (6) is insulated and the other surface isn't. Furthermore, one end of the specimen

(6) is fixed by the fixture (3) and the other end remains free. The free end of the specimen (6) has a predetermined distance from a vernier (51) of a vernier scale (5), which is mounted beside the testing cell (1) and touches the metal plate (4). A non-contact displacement sensor (16) is located close to the specimen (6) for sensing the displacement of the free end of the specimen (6) and generating a displacement signal during an electroplating process.

The apparatus employed allows analysis of the internal stress within the coating on one surface of the specimen subjected to electrodeposition and therefore may be used to help control the subsequent internal stress of electrodeposited materials. The equipment comprises a testing cell filled with the electroplating solution, and a mounting fixture for fixing one end of the specimen, leaving the other end free. A programmable power supply is then connected to the specimen and the cathode placed in the electroplating bath. The apparatus uses a non-contact linear sensor for measuring the displacement of the free end of the specimen when the current is applied. This sends a signal to the programmable power supply, and the displacement is used to calculate the internal stress within the coating. Two current channels of different current density alternately supply current to the specimen and a multi-layer deposit of two distinctly different alloys of different internal stress is produced. Therefore, it is possible to control the internal stress in the deposit by variations in the current densities used.

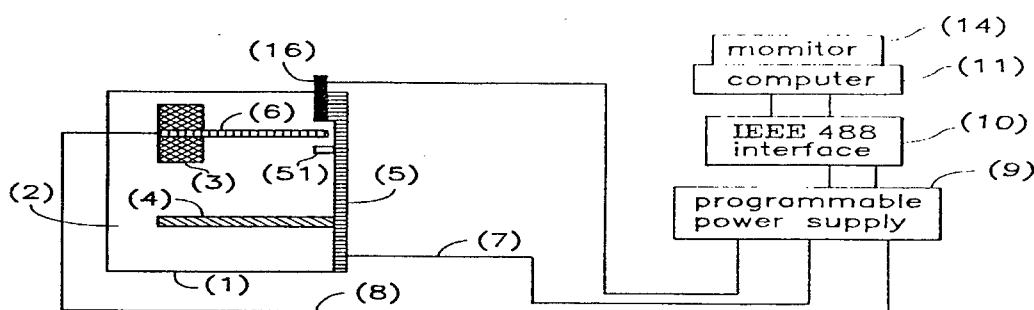


Fig. 1. Diagram of apparatus for analyzing the internal stress of a coating

The process of this study may be used to an advantage with conventional pulsed current, superimposed A.C. or periodic reverse plating. The process may also be advantageous for controlling the thickness of electrodeposits, since this property is also measured automatically and varies with current density. According to this study it has been found that automatic stress control in an electrodeposition may be achieved using the method and apparatus described herein, employing a single electroplating bath for electrodeposition of a multi-layer coating. Figure 1 illustrates the apparatus for measuring the internal stress of the electrodeposit.

A testing cell with a fixture for the specimen and a metal anode is used to electroplate the specimen. One end of the specimen is fixed and the other end remains free. The specimen is insulated on one side. A non-contact displacement sensor is located close to the specimen to monitor the movement at the free end of the specimen during an electroplating process. A computer programmable power supply is connected to a personal computer through an IEEE488 interface. The inductive linear sensor switches a relay to control two output channels to supply current to the substrate. During electroplating of the uninsulated surface of the specimen an internal stress is generated within the coating. The measuring sensor sends a displacement to the computer through the programmable power supply and the IEEE488 interface and the computer calculates the internal stress within the coating using the following formula:

Macro-stress, σ

$$\sigma = \frac{E_s T_s^2 + E_c T_s T_c \left(\frac{4 + 6T_c}{T_s + 3T_c^2 T_s^{-2}} \right)}{3L^2} \left(\frac{f}{T_c} \right) \frac{1}{(1 + 5T_c/3T_s)}$$

where

- S = Macro-stress
 E_s = Modulus of elasticity of substrate
 E_c = Modulus of elasticity of coating
 T_s = Thickness of coating
 T_c = Thickness of substrate
 L = Length of substrate
 f = deflection of free end of strip

The study will now be described in relation to the electrodeposition of a nickel-phosphorous alloy onto a copper substrate. The nickel phosphorous alloy plating solution was prepared by dissolving the following analytical grade compounds in distilled water in the amounts indicated in Table 1.

Table 1. Composition of Plating Solution

Compound	g/l
Nickel Sulfate	160
Nickel Chloride	50
Phosphoric Acid	50
Phosphorous Acid	60

The plating solution was then placed into a 150 mL electroplating cell. The cathode, a 0.25mm thick sheet of copper, insulated on one side, and formed in a U-shape as shown in Figure 2 was clamped at one end with the free end facing the inductive linear sensor. A rectangular nickel anode with approximate dimensions of 40mm high by 20mm wide was immersed in the bath at the far end of the plating cell. The two output channels, supplying current to the substrate, were set at 0.4A, giving an equivalent current density of 0.05A/cm², and at 1.4A, equivalent to a current density of 0.4A/cm². The inductive linear sensor was adjusted to two independent switching points corresponding to +0.05mm and -0.025mm displacements at the free end of the copper strip. By coupling the terminal output to the relay of the ON/OFF toggle switch a -0.025mm deflection sensed as one ON position and activated the high current channel. Similarly the low current channel was activated when the toggle switch at the other ON position sensed a +0.05mm deflection at the free end of the substrate. Hence a basic current pattern is achieved.



Fig. 2. Schematic diagram of the cathode.

RESULTS

- (1) The first charge burst of the low current of 0.4A produced a high phosphorous content ($P > 14\%$) nickel-phosphorous alloy with a compressive stress of about 7Kpsi (a). The second charge burst, of the high current of 1.4A produced a low phosphorous content ($P < 7\%$) nickel-phosphorous alloy with a tensile stress of about 15Kpsi (b). Repeating steps (a) and (b) produces a multi-layer deposit.
- (2) The high and low current densities for this process, according to (1) were found to be 0.4A/cm² and 0.05A/cm² respectively. The pulse train of Figure 3 consists of a repeating sequence of pairs of charge bursts of current pulses.

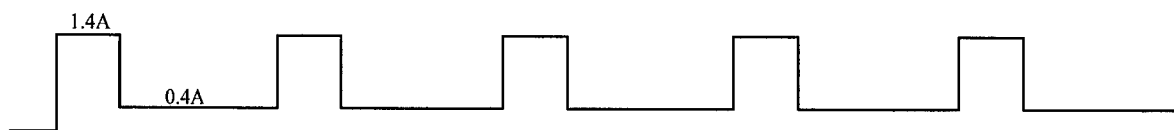


Fig. 3. Schematic Timing Diagram of Electrodeposition Current

- (3) A nickel-phosphorous plating film produced according to (1) comprised of two different kinds of films. The first film on the substrate had a high phosphorous content ($P > 12\%$) produced by each charge burst of low current and is compressively stressed. The second film had a lower phosphorous content ($P < 7\%$) deposited by each charge burst of higher electric current and is tensile stressed.
- (4) The multi-layer deposit produced according to (1) is essentially a sequence of repeating layers, each layer containing a high phosphorous content nickel-phosphorous alloy and a low phosphorous content nickel-phosphorous alloy. The figure 4 illustrates the cross-sectional view of layer deposit produced according to the process of Figure 3.
- (5) The pH of the process described in (1) was lower than a pH of 2 and was best adjusted by additions of phosphoric acid.

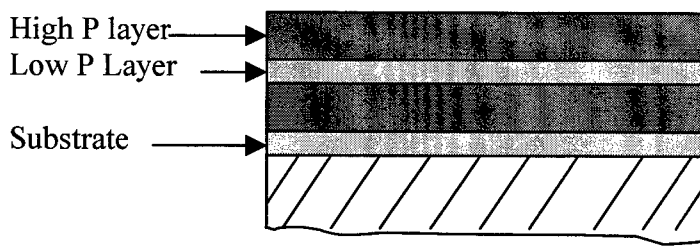


Fig. 4. Schematic cross-sectional view of layer deposit.

ACKNOWLEDGMENTS

The authors are grateful for the support of this research by the National Science Council, Republic of China, under Contract NSC83-0117-C-150-060E.

REFERENCES

1. G.G. Gawrilov, 1979. "Chemical (Electroless) Nickel Plating", Redhill, England, Poortcullis Press.
2. W.H. Metzger, 1959. "Characteristics of Deposit", American Society for Testing and Materials, Philadelphia.
3. G.D. Jarrett, 1966. "Electroless Nickel Plating", Industrial Finishing, 18, p.41.
4. L.G. Fitzgerald, 1960. "Nickel Plating Products", Finishing vol. 13, p.68.
5. W.H. Roberts, 1964. "Coating Beryllium with Electroless Nickel", U.S. Atomic Energy Commission, Rep478.
6. R.N. Duncan, 1981. "Properties and Applications of Electroless Nickel Deposits", Finisher Management, 26, p.5.

Electrolytic ZrO₂ Coating on Co-Cr-Mo Implant Alloys of Hip Prosthesis

S. K. Yen, M. J. Guo, and H. Z. Zan

Institute of Materials Engineering, National Chung Hsing University,
Taichung, 40227, Taiwan, R.O.C.

ABSTRACT

An electrolytic ZrO₂ gel has been coated on ASTM F-75 Co-Cr-Mo alloy specimens in 0.1 M ZrO(NO₃)₂ solution at pH 2.2 and a current density of 2 mA/cm². The electrolytic ZrO₂ gel coating was annealed at 623 – 973 K for 120 min. in air, then evaluated by electrochemical polarization in Hank's solution; wear tests with UHMWPE (ultra-high molecular-weight polyethylene) under a load stress of 50 Mpa; scratch tests; and morphology observations. The crystal structures of cobalt oxide and ZrO₂ were analyzed by XRD (X-ray diffraction). The ZrO₂-coated specimen annealed at 773 K for 120 min shows a good adhesion of 610 MPa on Co-Cr-Mo substrate, a lower wear loss of UHMWPE and a higher protection potential than the uncoated specimen in Hank's solution. A monoclinic structure with (111) preferred orientation parallel to the sheet plane was observed at 623 K ≤ T ≤ 673 K while a tetragonal structure of ZrO₂ was detected at T ≥ 773 K. For T ≥ 973 K, a more stable monoclinic structure with random orientation mixed with a tetragonal structure was observed.

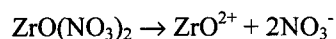
Keywords: electrolytic ZrO₂ coating, Co-Cr-Mo implant alloy, hip prosthesis

INTRODUCTION

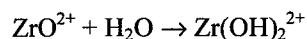
Cobalt-based alloys are widely used in total hip prostheses (THP), total knee replacements, dental implants, bone screws, staples, plates, and as support structures for heart valves because of their good wear and corrosion resistance, fatigue strength, and biocompatibility [1]. In THP surgery, an ultra-high-molecular-weight polyethylene (UHMWPE) socket is placed into the acetabulum, while a metallic hip stem consisting of a shaft and ball is placed in the femur. However, prostheses with bearing surfaces still release wear products into joint cavities. The gradual accumulation of these debris particles generates granulomatous reactions in the tissue surrounding the prostheses, and can cause loosening and osteolysis [2-5]. Variable amounts of plastic and metal are released due to wear and corrosion.

To cope with these problems, Boatin [6] in France developed a hip prosthesis with cup and ball made of alumina ceramic at the beginning of the 1970s. Semlitch et al. [7] determined that the wear rate of polyethylene (PE) against alumina ceramic is about 20 times lower than that of PE against Co-Cr-Mo alloys. They suggest that this favorable tribological behavior of ceramic in contact with PE may be due to better corrosion resistance, wettability with liquids, and scratch resistance of the ceramic materials compared to those of metallic implant materials. However, alumina ceramic exhibits a brittle tendency and is sensitive to microstructural flaws [8]. Currently, zirconia ceramic is being recognized for its high strength and surface finish, making this material potentially suitable for the highly loaded environments found in joint replacement [8-11]. A dip coating of zirconia gel exhibits a shear adhesion strength of 275 MPa on Ti-6Al-4V [12] and shows a better corrosion resistance on AISI 1006 and 304 stainless steel [13]. Besides, another competitive method, an electrolytic coating of zirconia gel on Ti-6Al-4V also shows a good adhesion of 580 MPa and a good corrosion resistance in 5 M HCl, 5 M H₂SO₄ or 0.6 M NaCl aqueous solution [14]. The electrolytic mechanism has been suggested [15] and modified [16] as follows:

(i) dissociation of zirconyl salt



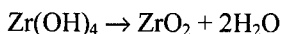
(ii) hydrolysis of the zirconyl ion



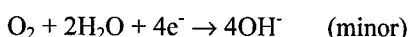
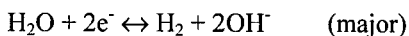
(iii) interaction with OH



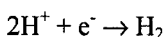
(iv) dehydration of the hydroxide



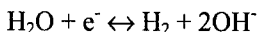
The main source of OH⁻ comes from the cathode:



And the best coating efficiency is found after the H⁺ diffusion limit occurs in



and the on set of the reduction of H₂O:



In this study, the characterization of electrolytic ZrO₂ coating on Co-Cr-Mo including crystal structure, corrosion resistance, and wear resistance is reported.

EXPERIMENTAL

Sample Preparation

An ASTM F-75 Co-Cr-Mo sheet, as received, was used as a substrate of the ZrO₂ electrolytic coating. Its chemical composition is given in Table I. The sheet thickness is 2.0 mm with a grain size about 5 μm ± 2 μm. The sheet was cut into discs with a diameter of 14 mm for corrosion tests, and 28 mm for wear and scratch tests. All specimens were polished to a mirror finish with 1 μm Al₂O₃ powder, then degreased by detergent and further ultrasonically cleaned in deionized water and acetone, then dried by an N₂ gas gun.

Table I. The nominal chemical composition of Co-Cr-Mo.

Element	Cr	Mo	Mn	Fe	C	Si	Mg	Ni	Co
Wt. %	26-30	5-7	1.0	0.75	0.35	1.0	1.0	0.25	Balance

Electrolytic Deposition and Annealing

The electrolytic deposition of ZrO₂ was conducted in a naturally aerated solution of 0.125 M ZrO(NO₃)₂, at pH 2.2 and a cathodic current density of 2 mA/cm² for 500 sec by using an EG&G M273A Potentiostat and M352 software. The alloy disc was the cathode, graphite was the anode, and saturated calomel was the reference electrode. The above electrolytic condition gave the most efficient deposition in our experiment. The specimens with Zr(OH)₄ gel coating were naturally dried in air and annealed in air at 623 K, 773 K, 873 K, 973 K, and 1073 K for 120 min, respectively.

Scratch and Wear Tests

Some specimens were tested by scratch (Teer St-2000) with a preload 2 N, load speed 50 N/min, scratch speed 20 mm/min, and end load 30 N. Reciprocating wear test was also conducted by Teer St-2000. UHMWPE was used as pin in diameter of 0.5 mm, ZrO₂ coated and uncoated specimens as disks with a contact surface area of 0.2 mm² with a load of 10 N (or axial stress of 50 MPa), at the sliding distance of 1 mm per cycle with 30 cycle/min. The drive speed and room temperature of 25 ± 1 °C were kept constant throughout the test. Frictional force between the UHMWPE specimen and the counterface was monitored by a strain gauge fixed on a leaf spring attached to the transverse bar holding the PE pin.

Corrosion Tests

All annealed specimens were potentiodynamically polarized with EG&G Model 273A M352 in aerated Hank's solution with compositions given as: NaCl 8.00 g/L, CaCl 0.14 g/L, KCl 0.40 g/L, NaHCO₃ 0.35 g/L, Glucose 1.00 g/L, MgCl₂·6H₂O 0.10 g/L, KH₂PO₄ 0.06 g/L, MgSO₄·7H₂O 0.06 g/L, and Na₂HPO₄ 0.06 g/L. The cyclic polarization test was from -0.80 to +0.80 V, then back to -0.70 V at a scanning rate of 5 mV/sec. The first oxidation-reduction equilibrium potential E_{01} was derived when current density equals zero during the applied voltage increased (forward cycle). If the oxidation is due to the corrosion of electrode, this potential is also named corrosion potential E_{corr} . The second oxidation-reduction equilibrium potential E_{02} was derived when the current density was returned to zero again during the applied voltage decrease (backward cycle). This potential was also related to protection potential E_{pp} , if the oxidation is due to the pitting of electrode.

SEM and XRD

The surface morphology of specimens after scratch and wear tested was observed by scanning electron microscopy/energy dispersive spectroscopy (SEM/EDS, JEOL JSM-5400 Japan). The crystal structure of ZrO₂ on Co-Cr-Mo substrate was analyzed by X-ray Diffractometry (XRD, MAC MO3X-HF Diffraction Japan), wavelength of Cu K α (1.5418 Å), 2 θ from 20° to 100°, at scanning rate of 0.5 °/min, voltage of 40 kV, and current of 30 mA.

RESULTS

The surface morphology of deposited specimen conducted at a cathodic current density of 2 mA/cm² for 500 sec is shown in Fig. 1 (a). Those of post-deposited specimen annealed at 773 and 873 K are shown in Figs. 1 (b) and (c), respectively. The average thickness of ZrO₂ film is estimated as 1 μ m by weight-gain measurements. Representative polarization curves of uncoated and coated specimen annealed at 673 K, 773 K, and 873 K for 120 min in Hank's solution are shown in Fig. 2. From those curves, the first redox potential E_1 (or corrosion potential E_{corr}), the second redox potential E_2 (or protective potential), exchange current i_0 (or corrosion current), cathodic polarization slope β_c , anodic polarization slope β_a are analyzed, as given in Table II. The surface morphology of specimen annealed at 873 K after polarization test is shown in Fig. 3.

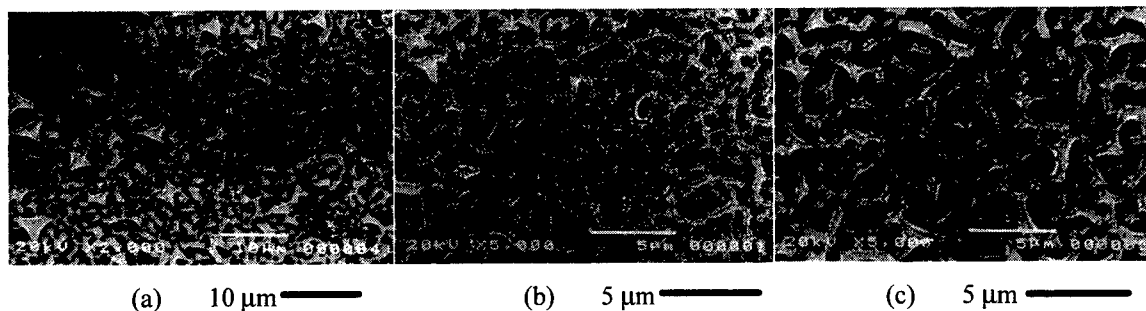


Fig. 1. SEM observations of deposited specimen (a) after natural drying, then annealed for 120 min at (b) 773 K and (c) 873 K.

Table 2. First redox potential E_1 , exchange current density i_0 , the second redox potential E_2 , cathodic polarization slope β_c , and anodic polarization slope β_a , derived from polarization tests in Hank's solution.

Annealed temperature	E_1 (E_{corr}) (V)	i_0 (i_{corr}) (μ A/cm ²)	E_2 (E_{pp}) (V)	β_c	β_a
As-received	-0.260	1.632	0.164	0.3596	0.4103
673 K	-0.316	4.448	0.192	0.3313	0.4568
773 K	-0.364	4.571	0.404	0.3261	0.5013
873 K	-0.380	4.706	0.156	0.2628	0.2892

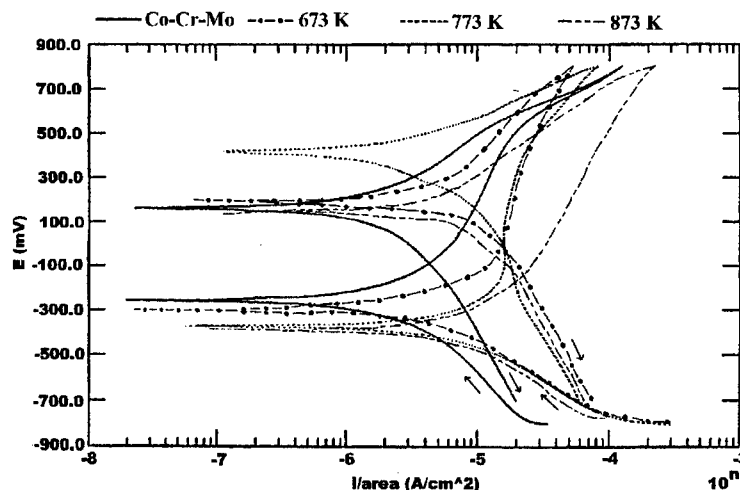


Fig. 2. Polarization curves of (a) uncoated specimen and ZrO_2 coated specimens annealed at (b) 673 K, (c) 773 K, and (d) 873 K tested in Hank's solution.

The surface morphology of ZrO_2 coated specimen annealed at 773 K after scratch tests are shown in Fig. 4. EDS mapping shows that a complete ZrO_2 film is still found at the end of scratch stage with a load of 30 N. The results of the reciprocating wear test are given in Table 3.

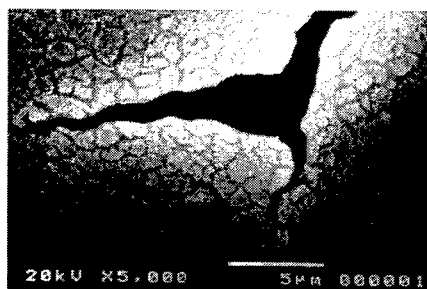


Fig. 3. SEM observations of specimen annealed at 873 K after polarization test.
(5 µm —)

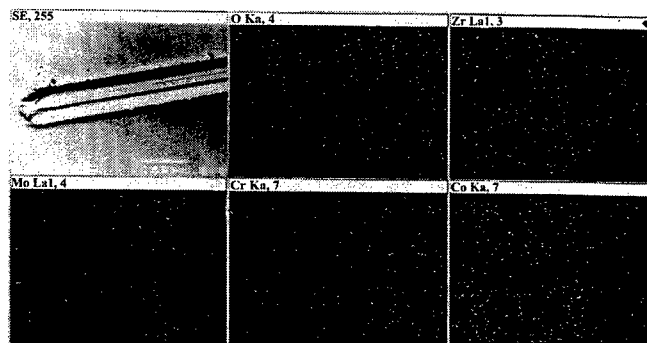


Fig. 4. SEM observations after scratch test of ZrO_2 coated specimens annealed at 773 K for 120 min.
(100 µm —)

Table 3. Results of the reciprocating wear test

	uncoated	ZrO_2 coated	uncoated	ZrO_2 coated	uncoated	ZrO_2 coated
Weight before test	0.07167	0.07602	0.10021	0.07559	0.08902	0.06451
Weight after test	0.07144	0.07587	0.09996	0.07548	0.08886	0.06444
Weight loss	0.00023	0.00015	0.00025	0.00011	0.00016	0.00007

The XRD of ZrO_2 on F-75 Co-Cr-Mo specimen annealed at 623, 673, 773, 873, and 973 K for 120 min are shown in Fig. 5 (1), (2), (3), (4), and (5), respectively. Only monoclinic ZrO_2 with $(\bar{1}11)$ preferred orientation parallel to sheet plane is found at $T \leq 873$ K and Co_2CrO_4 is also found at $T \geq 673$ K. Tetragonal ZrO_2 is found at $T \geq 773$ K. However, monoclinic ZrO_2 with random orientation mixed in tetragonal ZrO_2 was seen at $T \geq 973$ K.

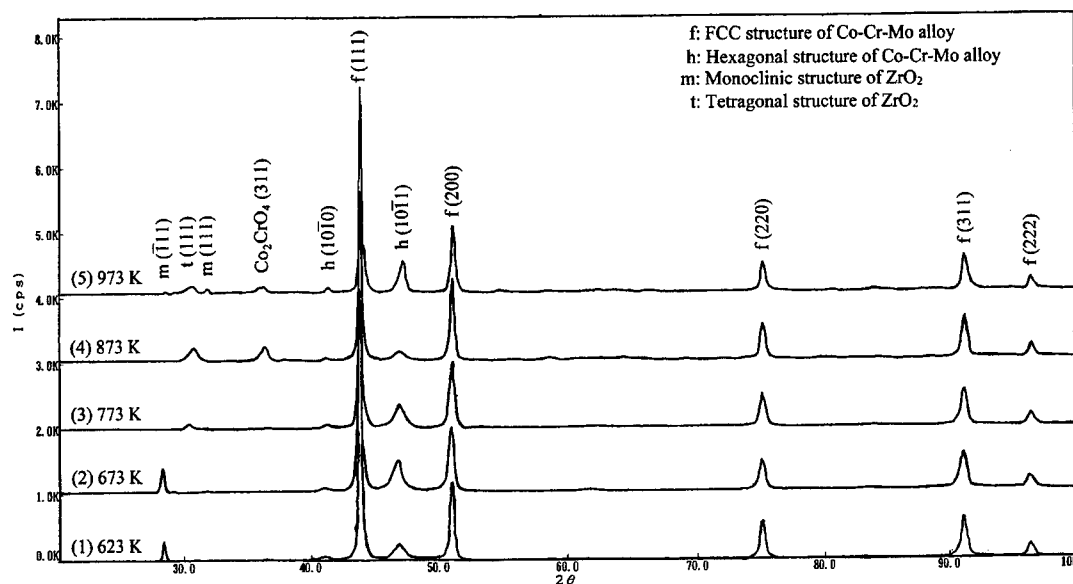


Fig. 5. XRD diagrams of ZrO_2 on Co-Cr-Mo substrate annealed at (1) 623 K, (2) 673 K, (3) 773 K, (4) 873 K, and (5) 973 K for 120 min.

DISCUSSION

Mud Crackings and Adhesion

Mud-crackings are found during natural drying in air, as shown in Fig. 1 (a). These crackings are possibly because the evaporation rate of H_2O from the surface of gel $R^e(\text{H}_2\text{O})$ is higher than the diffusion rate of H_2O from the bulk of the gel $R^d(\text{H}_2\text{O})$. If $R^e(\text{H}_2\text{O}) > R^d(\text{H}_2\text{O})$, a tension effect will be generated at the gel surface due to shrinkage of containing less H_2O than bulk gel layer and mud cracking will appear. If $R^e \leq R^d$, no tension effect will take place on the gel surface. Therefore, to control gel thickness (i.e. control R^d), control of the air humidity (i.e. control R^e) is a possible way to retard mud crackings. This argument is supported by an experiment in which a specimen with coated film thickness under $0.3 \mu\text{m}$, or another experiment in which the humidity was controlled at 95% showed no mud crackings. Another possible reason for mud cracking is the accommodating stress between the coating and substrate. The extent of mud crackings is changed slightly after annealing, as shown in Fig. 1 (a) as deposited, (b) annealed at 773 K, (c) annealed at 873 K. This means that the difference in thermal expansion between the metal substrate and ceramic coating will alter the extent of cracking. However, the number of cracks increased at 873 K, as shown in Fig. 1 (c). This is because oxidation of Co-Cr-Mo at 873 K is faster and the Co_2CrO_4 oxide will form at the mud-cracking due to a lack of ZrO_2 protection. The XRD diagram also indicates the increasing intensity of the Co_2CrO_4 peak, as shown in Fig. 5 (4) and (5). So Co_2CrO_4 will grow along the interface between the Co-Cr-Mo alloy substrate and the ZrO_2 coating, and break the ZrO_2 coating into even more pieces. On the contrary, most ZrO_2 is still found at the end stage of the scratch test, as shown in Fig. 4 (b). With a load of 30 N and $70 \mu\text{m}$ width of the scratch, the stress at this stage is about 610 MPa which is 17 times the loading stress on the hip joint during gait ($< 35 \text{ MPa}$).

Corrosion Resistance and Wear Loss

From Table 2, all post-deposited ZrO_2 specimens have shown a lower E_{01} and a lower cathodic polarization slope β_c . Corrosion potential E_{corr} is higher than the redox potential for $2\text{H}_2\text{O} + 2e^- \rightarrow \text{H}_2 + 2\text{OH}^-$ where $E_{\text{H}_2}^0 = -0.635$ (vs. AgCl, pH 7), as given in Table 2. A lower β_c means that reduction of H_2 is easier on the ZrO_2 coated specimen. No pitting potential was found in the cyclic polarization test (from -0.8 V to $+0.8 \text{ V}$, then back to -0.7 V) for ZrO_2 coated specimen annealed at 673 K and 773 K. Consequently, the protection potential E_{pp} means the second redox potential E_{02} is possibly a redox for $2\text{H}_2\text{O} + \text{O}_2 + 4e^- \leftrightarrow 4\text{OH}^-$ $E_{\text{O}_2}^0 = 0.598$ (vs. AgCl, pH = 3). On the other hand, a lower E_{02} was found on ZrO_2 coated specimen annealed at

873 K and uncoated specimen, and a higher current density was found at applied voltage ≥ 700 mV. Corrosion attack was observed on uncoated and ZrO_2 coated specimen annealed at 873 K, as shown in Fig. 3 but none was found on ZrO_2 coated specimen annealed at $T \leq 773$ K. No better corrosion resistance of coated specimen annealed at $T \geq 873$ K is likely due to the formation of Co_2CrO_4 which destroys the adhesion of ZrO_2 to the Co-Cr-Mo substrate, as shown in Fig. 1 (c) and Fig. 5 (4). From Table 3, the wear loss of UHMWPE for ZrO_2 -coated specimens is obviously less than that of the uncoated specimen since the friction coefficient of UHMWPE (0.13) to ZrO_2 -coated specimen is smaller than the (0.17) uncoated specimen.

Crystal Structures

XRD diagrams of Fig. 5 (1), (2), (3), (4), and (5) show ZrO_2 phase transformations among monoclinic with ($\bar{1}11$) preferred orientation, tetragonal, and random monoclinic crystal structures. At a lower annealing temperature $T \leq 673$ K, amorphous ZrO_2 was transformed monoclinic with ($\bar{1}11$) preferred, as shown in Fig. 5 (1) and (2). The tetragonal structure was detected at $T \geq 773$ K, as shown in Fig. 5 (3) and (4). Then the tetragonal was gradually transformed into a monoclinic with random orientation at $T \geq 973$ K, as shown in Fig. 5 (5).

Many arguments have been suggested to indicate ZrO_2 among amorphous, monoclinic and tetragonal crystal structures, such as particle size effect [17,18], precursor amorphous phase [19,20], anionic vacancy [21], and pH values [22-25], but there is no convincing one which can completely describe these results of this study. This is probably due to the different processing of ZrO_2 . Obviously, the structure of ZrO_2 in this study is dependent on the annealing time and temperature. This argument has been suggested before [24-26], but it has not been found before that a tetragonal structure shows metastable at $698 \text{ K} < T \leq 973 \text{ K}$, a temperature range below which a monoclinic structure is found. Possibly, the substrate effect which will affect the interfacial energy and then favor the nucleation of monoclinic ($\bar{1}11$) preferred orientation at $T \leq 698 \text{ K}$ and tetragonal ZrO_2 at $698 \text{ K} < T \leq 973 \text{ K}$ on it, should be also considered. In other words, the nucleation activation energy of a monoclinic with ($\bar{1}11$) preferred orientation or tetragonal structure is lower than that of a monoclinic with random orientation. Similar results were also found on ZrO_2 coated Ti-6Al-4V alloy [14]. Therefore, at lower annealing temperature and/or for a shorter time, a kinetic will dominate the phase transformation. However, at higher annealing temperature and for a longer annealing time, thermodynamics of a phase with much lower free energy will dominate the phase transformation, such as monoclinic structure with random orientation.

SUMMARY AND CONCLUSIONS

A new electrolytic coating method of ZrO_2 has been applied on an ASTM F-75 Co-Cr-Mo alloy to investigate its characteristics. Through the electrolytic coating, drying, annealing, polarization tests, surface observations, XRD analysis, and scratch tests, several conclusions are drawn:

1. The mud cracks are found during natural drying in air. To avoid these cracks, some moisture-controlled system must be applied or the thickness of the gel controlled to improve these problems.
2. Porous Co_2CrO_4 particles which are not protective and without good adhesion to the substrate was obviously found at annealed $T \geq 873 \text{ K}$, a vacuum annealing furnace should be helpful to avoid the serious oxidation on the Co-Cr-Mo substrate.
3. However, ZrO_2 post-coated Co-Cr-Mo specimen annealed at 773 K for 120 min has shown a higher protective potential for the polarization test in Hank's solution, less wear loss of UHMWPE, and a lower friction coefficient from 0.17 to 0.13.
4. For 120 min annealing, a monoclinic structure with ($\bar{1}11$) preferred orientation parallel to the Co-Cr-Mo sheet plane was detected only at $623 \text{ K} \leq T \leq 673 \text{ K}$, a tetragonal structure was found at $T > 773 \text{ K}$, and a more stable monoclinic structure with random orientation was detected with a mixed tetragonal structure at $T \geq 973 \text{ K}$.

ACKNOWLEDGEMENTS

The authors are grateful for the support of this research by National Science Council, Republic of China under contract No. NSC 87-2213-E-005-022.

REFERENCES

1. Aziz I. Asphahani, Haynes International, Inc., 1988. ASM Metal Handbook 9th ed., 13, 658-668.
2. U. E. Pazzaglia, L. Ceciliani, M. J. Wilkinson, and C. Dell'Orbo, 1985. Arch. Orthop. Traumat. Surg., 104, 164-174.
3. A. Pizzoferrato, L. Savarino, S. Stea, and C. Tarabusi, 1988. Biomaterials, 9, 314-318.
4. F. Betts, T. Wright, E. A. Salvati, A. Boskey, and M. Bansal, 1992. Clinical Orthopaedics and Related Research, 276, 75-82.
5. F. F. Henning, H. J. Raithel, K. H. Schaller, and J. R. Döhler, 1992. J. Trace Elem. Electrolytes Health Dis., 6, 239-243.
6. P. Boutin, 1971. Press Med., 79, 639.
7. M. Semlitsch, M. Lehmann, and H. Weber, 1977. J. Biomed. Mater. Res., 11, 537-552.
8. J. P. Torre, 1986. Proc. Materials Eng. Conf., 5-7th November 1985, London, UK, Mechanical Engineering Publication Ltd., Bury St., Edmund, Suffolk, 1986.
9. K. Shimizu, P. Kumar, M. Oka, Y. Kotoura, Y. Nakayama, T. Yamamuro, T. Yanagida, and K. Makinouchi, 1988. Transaction of the 3rd. World Biomaterial Congress, 21-25th. April, 1988, Kyoto, Japan, Vol. 11, p. 406.
10. K. Makinouchi, T. Yanagida, K. Shimizu, P. Kumar, Y. Kotoura, and M. Oka, 1984. Proc. Jap. Soc. Orthop. Ceramic Implants, H. Oonishi and Y. Ooi (eds.), 4.
11. P. Christel, A. Meunier, M. Heller, J. P. Torre, and C. N. Peille, 1989. J. Biomed. Mater. Res., 23, 45-61.
12. M. J. Filiaggi, R. M. Pilliar, D. Abdulla, 1996. J. Biomed. Mater. Res. (Appl. Biomater.) 33, 239-256.
13. F. Derdomo, P. de Lima, M. A. Aegerter, and L. A. Avaca, 1996. 13th Int. Corrosion Congress, Melbourne, Australia, 1996, paper 082.
14. S. K. Yen, T. Y. Huang, 1998. Materials Chemistry and Physics, 56, 214-221.
15. L. Gal-Or, I. Silberman, and R. Chaim, 1991. J. Electrochem. Soc., 138, 1939-1942.
16. S. K. Yen, 1997. 192nd Meeting Abstracts of The Electrochemical Society, 97-2, 636.
17. R. C. Garvie, 1985. J. Phys. Chem, 82, 218.
18. R. C. Garvie and M. V. Swain, 1985. J. Mater. Sci., 20, 1193.
19. J. Livage, K. Doi, and C. Mazieres, 1968. J. Am. Ceram. Soc., 51, 349.
20. E. Tain, M. Yoshimura, and S. Somiya, 1983. J. Am. Ceram. Soc., 66, 11.
21. M. I. Osendi, J. S. Moya, C. J. Serna, and J. Soria, 1985. J. Am. Ceram. Soc., 68, 135.
22. R. Srinivasan, R. D. Angelis, and B. H. Davis, 1986. J. Mater. Res., 1(4), 583.
23. B. H. Davis, 1984. J. Am. Ceram. Soc., 67, 168.
24. R. Srinivasan, M. B. Harris, S. F. Simpson, and B. H. Davis, 1988. J. Mater. Res., 3(4), 787.
25. S. S. Jada and N. G. Peletis, 1989. J. Mater. Sci. Lett., 8, 243.
26. P. Singh and S. K. Date, 1987. J. Mater. Sci. Lett., 6, 621.

A New Process to Produce Advanced Zirconia-based Ceramic Composites from Low-Value Minerals

S. M. B. Veiga*, M. M. Veiga**, A.C.D. Chaklader***, J. C. Bressiani*

*Instituto de Pesquisas Energéticas e Nucleares, IPEN/CNEN, São Paulo, SP, Brasil.

**University of British Columbia, Department of Mining and Mineral Process Engineering, Vancouver, BC, Canada

***University of British Columbia, Department of Metals and Materials Engineering, Vancouver, BC, Canada.

ABSTRACT

Knowledge about the relationships between microstructure and properties is important to develop structural ceramics. However, the type of processing and purity of the powder affect microstructure and consequently mechanical properties. In fact, composite ceramic powders have been developed to enhance quality of structural ceramics. Zirconia-based ceramics have shown significant improvement of fracture toughness. Al_2O_3 -SiC-ZrO₂ ceramic-ceramic composite has been prepared mixing powder of all three ceramic components. Moreover, this method results segregated mixtures, utilizes expensive pure powders and introduces health problem as SiC-whiskers are carcinogenic.

A method to produce Al_2O_3 -SiC-ZrO₂ powder composite by carbothermal reaction was investigated. Carbothermal reaction has been a creative technique to produce alumina-silicon carbide composite powder from inexpensive precursor materials such as kaolinite, kyanite, pyrophyllite, etc. The products obtained from carbothermal reactions have shown nanometric particle sizes, homogeneous mixture and most impurities were eliminated by volatilization. Zircon (ZrSiO_4), as an inexpensive source of zirconia, was mixed with kaolinite-carbon or kyanite-carbon to produce zirconia-based composites. Unfortunately zirconia cannot be obtained directly from carbothermal reaction of these minerals as the reaction to produce zirconium carbide is favored. Instead, this new process obtains Al_2O_3 -SiC-ZrC composite powder at temperatures above 1500°C at 1 atm. However, a subsequent controlled oxidation step can transform ZrC of this powder into a mixture of monoclinic and tetragonal ZrO₂. Thermodynamic data were generated to support test results.

The Al_2O_3 -SiC-ZrO₂ powder with 7.9%vol ZrO₂ and 23.4%vol SiC was sintered by hot pressing at 1800 °C resulting in pellets with 30% higher fracture toughness than the ones made of Al_2O_3 -SiC composite. This encouraging result led to conclude that carbothermal reaction is a significant process to obtain ceramic composites by using different types of inexpensive minerals.

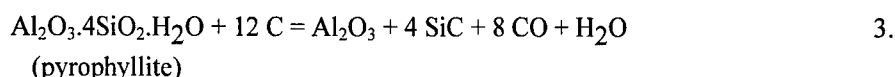
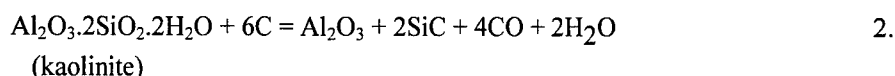
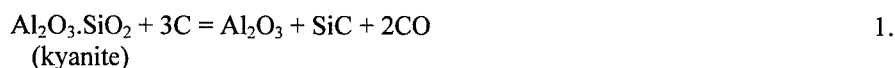
INTRODUCTION

Developments in advanced ceramics have resulted in more durable materials and better mechanical properties than the traditional materials. In fact, advanced ceramics are expanding their markets competing with metals in many structural applications where resistance to corrosion and performance at high temperatures are required properties. The American market of structural ceramics is estimated to be around US\$ 1.5 billion and rising. Automobile engines, cutting tool bits, pump and valve components, seals, have been the main focus of the advanced ceramic industries.

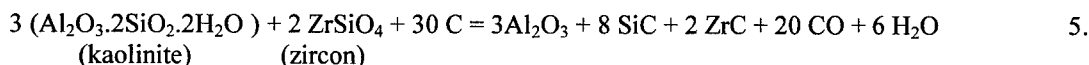
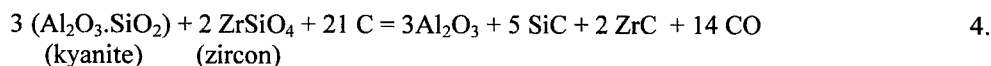
Two main impediments for a larger use of advanced ceramics in structural applications are the brittleness and relative high price of these materials. Polycrystalline ceramics, including oxides and covalent ceramics, have fracture toughness, usually below $6 \text{ Mpa}^{1/2}$ which is very low when compared with metals, 15 to 150 $\text{Mpa}^{1/2}$. The fragile nature of ceramics has led researchers to develop alternative ways to improve toughness such as reinforcement with fibers, whiskers or dispersion of a second phase. Mixtures of two or more different types of ceramic powders with different morphology have created stronger composites. Ceramic-

ceramic composites have shown considerable improvement in the strength and fracture toughness with introducing of zirconia. Such studies have led to the belief that ceramic-ceramic composites are the best materials for high temperature structural applications. One group of such composites are based on alumina (Al_2O_3) and silicon carbide (SiC) whiskers, which are very expensive ($\sim \$200/\text{kg}$), not easily available and also considered to be carcinogenic.

Chaklader et al. [1] have shown that advanced ceramic composites, based on alumina and SiC , can be synthesized by carbothermal reactions from cheap precursor materials such as kyanite, kaolinite, pyrophyllite, etc. The final reactions that express the process are:



It is well known that zirconia (ZrO_2), in particular in the tetragonal form, can improve ceramic toughness due to the martensitic transformation (tetragonal to monoclinic zirconia) that increases phase volume and creates sites that are resistant to crack propagation [2,3,4]. In an initial approach, it was thought that zircon (ZrSiO_4), which is a common Zr mineral, could be introduced in the carbothermal reaction to form zirconia in the composite matrix. However, preliminary tests have shown that Al_2O_3 - SiC - ZrO_2 cannot be obtained directly from carbon, kaolinite or kyanite and zircon reaction. Formation of Al_2O_3 - SiC - ZrC is always resulted in an inert atmosphere. A further oxidation step is necessary to transform ZrC into ZrO_2 [5,6]. The reactions involved in the carbothermal process are derived and expressed by the following equations:



According to equations 4 and 5, the products have the following composition: 43% and 37% Al_2O_3 ; 28% and 38% SiC ; 29% and 25% ZrC , respectively.

Volatile silicon monoxide is an intermediate phase for SiC formation. Reaction of this gaseous phase with carbon produces silicon carbide according to the equation (6). The availability of SiO in the system is a parameter difficult to control. Part of forming carbon monoxide reduces more silica and part is lost in the fumes [7].



This work presents results of an experimental investigation to synthesize Al_2O_3 - SiC - ZrO_2 composite from kaolinite or kyanite and zircon and a brief discussion about the relevant thermodynamic concepts related to carbothermal reactions.

EXPERIMENTAL PROGRAM

Mixtures of zircon, kaolinite or kyanite (Table 1) and lamp black carbon were ground in a porcelain ball mill resulting in a homogeneous powder as fine as $1\text{ }\mu\text{m}$. Around 32 g of the mixture were used to prepare pellets with diameter of 50.7 mm. Pellets were placed in graphite crucibles and heated in an induction vertical furnace with Argon atmosphere. Temperatures used for synthesis ranged from $1460\text{ }^\circ\text{C}$ to $1600\text{ }^\circ\text{C}$ with different residence time (120, 150 and 180 min) in the induction furnace.

Table 1. Chemical composition of kaolinite, kyanite and zircon.

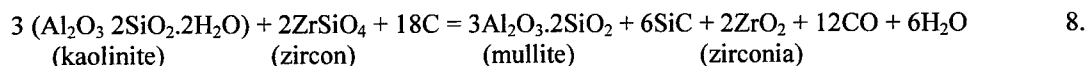
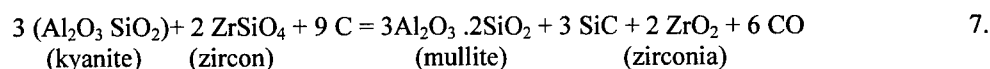
Substance	Kaolinite (wt%)	Kyanite (wt%)	Zircon (wt%)
Al ₂ O ₃	38.74	57.56	0.01
SiO ₂	44.54	39.87	33.76
CaO	0.19	0.11	0.03
Fe ₂ O ₃	0.52	0.75	0.01
MgO	0.05	0.15	0.01
K ₂ O	0.16	0.05	0.01
Na ₂ O	0.05	0.07	0.01
P ₂ O ₅	0.12	0.17	0.11
TiO ₂	1.63	1.26	0.11
MnO	-	-	0.03
Ba	-	-	0.01
ZrO ₂	-	-	65.23
H ₂ O	14.0	-	-

RESULTS AND DISCUSSION

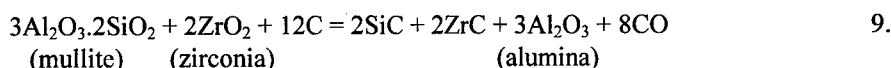
Carbothermal Reactions

The firing conditions have varied to determine the reaction steps to form Al₂O₃, SiC and ZrC composite. The effect of changing composition of the reagents (carbon, kaolinite or kyanite and zircon) was tested in 56 experiments and no other phase was produced, but rather, the same products with different compositions of Al₂O₃, SiC and ZrC. It was noticed that zircon is transformed into zirconia at the same temperature level (around 1300 °C) as when mullite and free silica are formed from kaolinite or kyanite. ZrC is not formed at temperatures below 1300 °C. At temperatures slightly below 1500°C and residence time between 120 and 180 min., formation of alumina, SiC, ZrC, zirconia and mullite was always observed.

Transformation of kaolinite or kyanite into mullite (3Al₂O₃.2SiO₂), free silica (SiO₂) and subsequently gaseous SiO formation is assumed as the first reaction step occurring at temperatures higher than 1300 °C [8]. At this temperature level, gaseous SiO reacts with carbon producing silicon carbide (SiC). Zircon is also transformed into zirconia at the same temperature range as observed by the presence of mullite and zirconia. Based on these experimental observations, the following intermediate reactions are proposed:



At temperatures higher than 1500 °C, the products from the intermediate reactions (7 and 8) react with carbon to produce α-alumina (Al₂O₃), zirconium carbide (ZrC) and additional β-silicon carbide (SiC). Zirconium carbide is formed as a product of zirconia reduction :



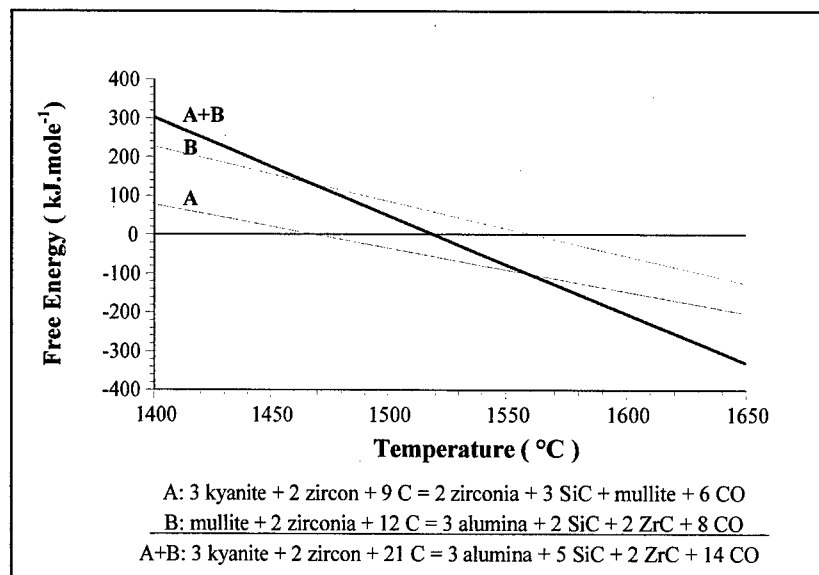


Fig. 1. Free energy diagram of carbothermal reactions between kyanite and zircon

Figures 1 and 2 are derived from thermodynamic data of intermediate reactions. At temperatures lower than 1560 °C, mullite and zirconia are predominantly formed. Alumina and ZrC are stable at temperature levels higher than 1560 °C. Partial CO pressure was not measured inside the induction furnace during the synthesis process, but it plays significant role on the carbothermal reactions. For example, assuming the partial pressure as 1 atm, at 1500 °C, the free energy of reaction A+B, in Fig.2 (which is the equation 5), is -1205 kJ.mole⁻¹. However if P_{CO} is considered hypothetically as 10⁻² atm, the free energy of the resulting reaction decreases to -2570 kJ.mole⁻¹ indicating that the reaction becomes more favorable. Actually the partial pressure of CO is definitely lower than 1 atm as Argon is injected into the synthesis chamber. The P_{CO} controls the conditions of ZrC formation. This can be a plausible explanation for the formation of ZrC at temperatures below 1560 °C as observed previously [9].

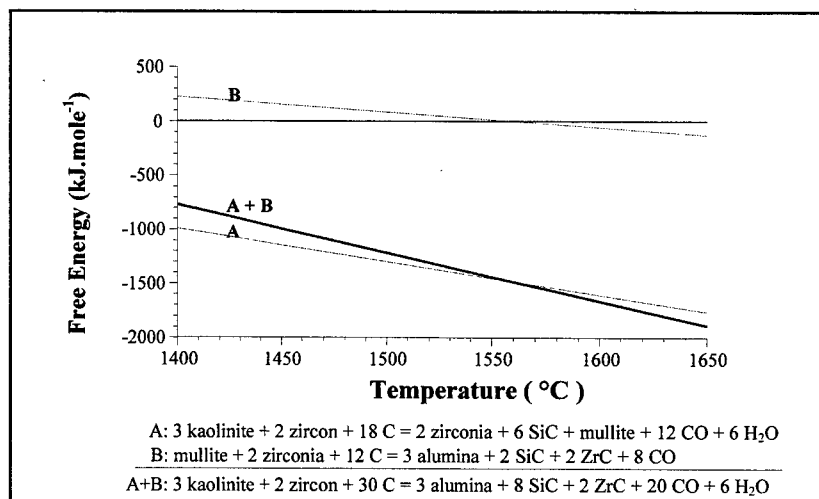


Fig. 2. Free energy diagram of carbothermal reactions between kaolinite and zircon

The equilibrium ($\Delta G^\circ = 0$) temperature of the final equation (A+B, in Fig. 1) of carbothermal reaction of kyanite and zircon is calculated as 1520 °C. This indicates that the reaction becomes spontaneous above this temperature. However, until 1560 °C, the free energy of the intermediate reaction A (formation of zirconia, SiC and mullite) is lower than the reaction A+B. Thus, only above 1560 °C, does reaction A+B, becomes thermodynamically favorable.

Reaction yields were calculated by dividing the amount of SiC formed in the system analyzed by quantitative x-ray diffractometry by the theoretical %wt SiC according to equation 4: 31.2%SiC and equation 5: 37.8% SiC. The yields ranged from 85% to 99%. Losses of SiO (gaseous) as well as the formation of amorphous SiO₂ are the main factors controlling reaction yield.

Oxidation Tests

Oxidation of ZrC was studied in both an TGA equipment and electric furnace using analytical grade ZrC. Pure β -SiC (grain size of 0.5 μ m) was also used to evaluate the degree of oxidation of this compound under oxidizing conditions (air and oxygen). Using TGA equipment, no relevant oxidation was observed when SiC particles were heated up to 900 °C with oxygen (Figure 3). In contrast, 26 mg of pure ZrC increased 20% of weight when oxidized at temperatures higher than 600 °C with oxygen forming monoclinic and tetragonal ZrO₂.

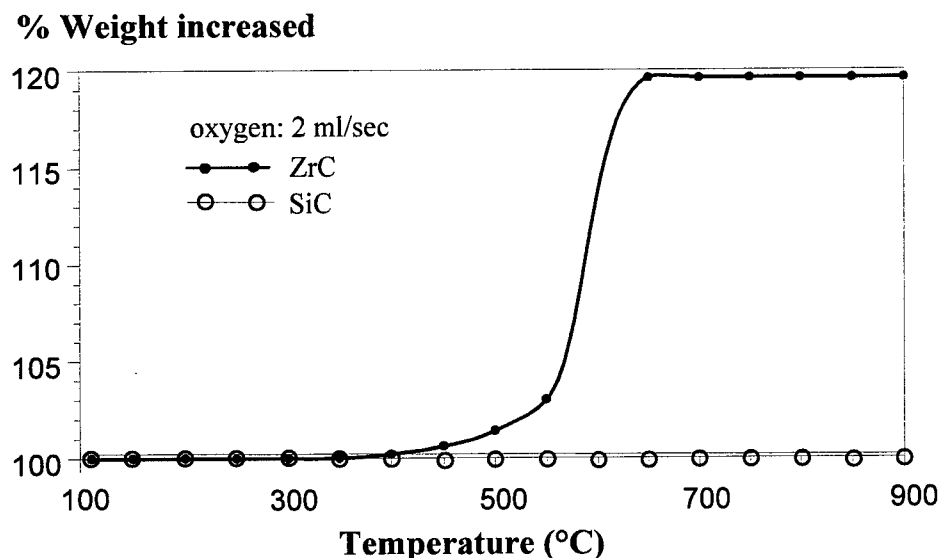


Fig. 3. Oxidation of ZrC and SiC Thermogravimetric Analysis.

In an electric furnace at 700 °C, it was observed by x-ray diffraction and weight gain control that, 20 g of a thin layer of ZrC spread on an alumina tube were partially oxidized after 12 hours of heating with oxygen injection. Monoclinic zirconia was resulted together with an intermediate compound (supposedly zirconium oxi-carbide) which is likely formed as a thin layer onto ZrC particles that hinders oxidation progress. This phase, as previously observed by other authors [5, 10](Bartlett et al., 1963; Shimada and Ishii, 1990), is formed as a result of poor oxygen diffusion into ZrC particles. Complete oxidation of ZrC was only achieved introducing 20% of a strong oxidizing agent (KNO₃) in the system. In this case, the oxidation time was reduced to 5 hours, temperature to 500 °C and a mixture of tetragonal and monoclinic zirconia was obtained. Once oxidation step was completed, the excess oxidizing agent was easily removed from the final composite powder with hot water. No relevant oxidation was observed when pure β -SiC was oxidized under the same conditions as zirconium carbide.

X-ray diffraction pattern of the composite powders oxidized with potassium nitrate, clearly indicated formation of both monoclinic and tetragonal zirconia (Figure 4). The concentration of the tetragonal phase in different powders varied from 30% to almost 100% of the total amount of zirconia. Even with the oxidizing agent, the more dispersed was the powder in the oxidation process, the more tetragonal zirconia was obtained (Figure 5).

Particle Analysis

It has been earlier reported [11] that the particles of SiC and Al₂O₃ formed from reduction of aluminosilicates are normally very small, in the range of below 1 μ m. This is specially true if lamp black is used as the carbon source. Furthermore, it was also observed that the particle size and shape of SiC can be controlled by using morphologically controlled carbon powder in the carbothermal reduction process [8].

In this investigation it is also considered that particle size of the final phases, specially that of ZrC and ZrO_2 , may have important bearings in developing tough ceramic composite materials. To enhance the fracture toughness by transformation toughening, the size and size-distribution of the ZrO_2 phase in the matrix are critical parameters. With this in view, particle size analyses were carried out on specimens containing ZrC and also on oxidized samples containing only ZrO_2 (in addition to SiC and Al_2O_3 being present in all systems). It is also well known that very fine particles of zirconia ($< 0.2 \mu\text{m}$) are normally in the metastable state of tetragonal form at ambient temperature.

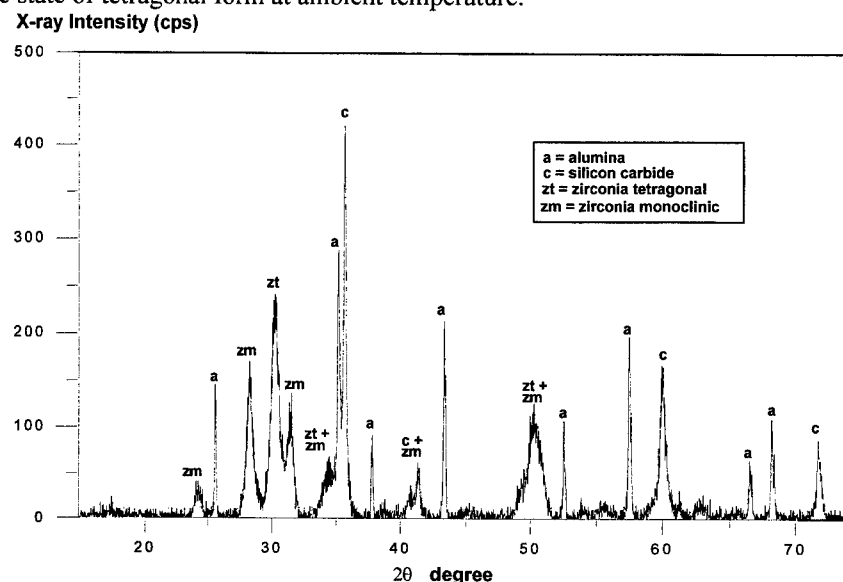


Fig. 4. X-ray pattern showing the presence of both tetragonal and monoclinic zirconia in the system alumina-silicon carbide-zirconia (after oxidation)

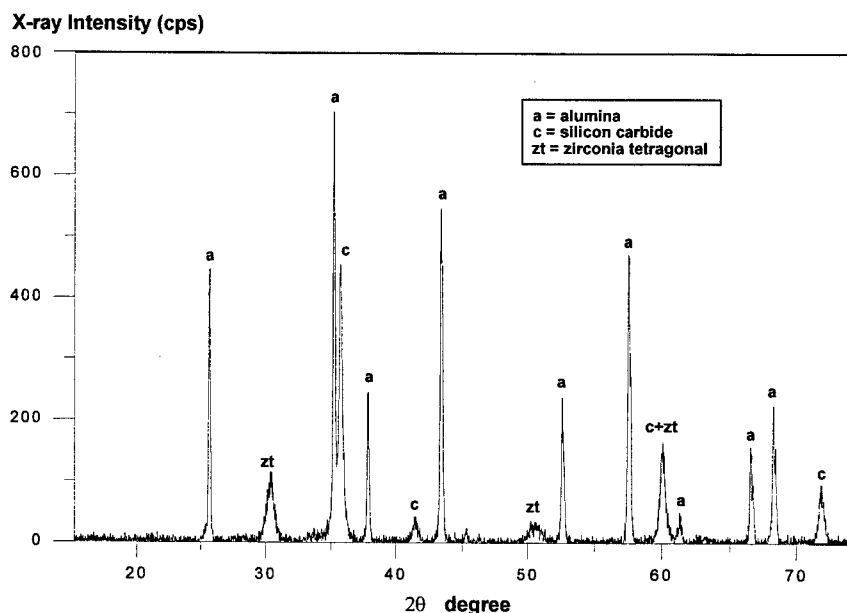


Fig. 5. X-ray pattern showing the predominance of tetragonal zirconia in the system alumina-silicon carbide-zirconia (after oxidation).

Figure 5 shows the x-ray diffraction pattern of a sample in which all ZrO_2 is in the form of tetragonal phase. Particle size analyses carried out on the powders containing Al_2O_3 , SiC and ZrC (before oxidation) have shown that 50% of particles before oxidation are smaller than $1 \mu\text{m}$. The larger particles are identified as Al_2O_3 by SEM coupled with EDS. As a matter of fact, SEM photographs with EDS analysis showed that

almost all ZrO_2 particles formed after oxidation are finer than $0.5 \mu\text{m}$. This suggests finer particles of ZrC transform preferentially into tetragonal zirconia [12]. Not infrequently, ZrC particles form agglomerates in the oxidation step. This favors formation of monoclinic ZrO_2 .

Fracture Toughness of Pellets

The densification of Al_2O_3 -SiC-ZrO₂ powder utilized a hot pressing process at 1800°C and 20 MPa to produce pellets with diameter of 20 mm. About 33 pressing tests were conducted. It was observed that a significant part of the tetragonal zirconia was transformed into monoclinic phase after hot pressing. Adding 5% mole of yttria (Y_2O_3) related to zirconia moles, most zirconia in the pressed powder was obtained in the tetragonal form. The final pellet achieved 97% densification.

Fracture toughness was evaluated calculating the K_{IC} , critical stress intensity factor, by using Vickers indentation method. Three composite powders were synthesized to prepare yttria-stabilized pellets of powders with 0, 7.9% and 17.1% vol ZrO_2 and 54.6%, 23.4% and 43.6% vol SiC. The K_{IC} results shown in Table 2 are in agreement with previous works, that excess zirconia does not contribute to increase fracture toughness. According to Wang and Stevens [13], the fracture toughness of Al_2O_3 -ZrO₂ composites increases linearly for mixtures up to 10% vol ZrO_2 and then decreases for mixtures with additional zirconia. Deleterious effects on fracture toughness of Al_2O_3 -SiC-ZrO₂ composites with more than 25% vol of SiC and 15% vol of ZrO_2 have also been observed [14]. High contents of ZrO_2 and SiC derive conjugated effects on the micro-crack propagation mechanisms that, ultimately, reduce toughness.

Table 2. Fracture toughness of ceramic composites with and without zirconia

Composite	ZrO ₂ (%vol)	SiC (%vol)	Original mixture of reagents in the carbothermal reaction	K_{IC} (MPa.m ^{1/2})
Al_2O_3 -SiC	0	54.6	78.2% kaolinite + 21.8% C	5.8 ± 0.3
Al_2O_3 -SiC-ZrO ₂	7.9	23.4	70.2% kaolinite + 8.5% zircon + 21.3% C	7.5 ± 1.1
Al_2O_3 -SiC-ZrO ₂ *	17.1	43.6	51.6% kaolinite + 24.4% zircon + 24.7% C	5.5 ± 0.8

* stoichiometric composition according to equation (5)

CONCLUSIONS

Al_2O_3 -SiC-ZrC is preferentially formed by carbothermal reaction of kaolinite or kyanite and zircon in Argon atmosphere. An oxidation step, using a strong oxidizing agent, such as KNO_3 , is needed to convert effectively ZrC into ZrO_2 to obtain powder of Al_2O_3 -SiC-ZrO₂ as final product. Without potassium nitrate, a passive thin layer is likely formed onto ZrC particles hindering oxidation. It is inferred by SEM-EDS observations together with x-ray diffraction that tetragonal zirconia is formed preferentially from small particles of ZrC. Furthermore, in the oxidation process, monoclinic zirconia is predominantly formed from larger particles of ZrC or from agglomerates. Particle dispersion in the oxidation process is a key step in obtaining tetragonal zirconia in the composite powder.

Thermodynamic data have supported the test results, in which ZrO_2 and mullite are formed at temperatures below 1560°C in the carbothermal synthesis. The Al_2O_3 -SiC-ZrC composite powder is preferentially formed at temperatures above 1560°C at 1 atm. Variations in the amounts of kaolinite and zircon did not form different products but always the same composite (Al_2O_3 -SiC-ZrC) powder with different compositions.

The fracture toughness of pellets of Al_2O_3 -SiC-ZrO₂ with 7.9%vol ZrO_2 and 23.4%vol SiC densified by hot pressing increased about 30% compared with Al_2O_3 -SiC composites also obtained by carbothermal process.

ACKNOWLEDGMENT

The authors acknowledge financial support of the Brazilian National Research Council (CNPq) for part of the experimental work.

REFERENCES

1. Chaklader, A.C.D.; Gupta, S.D.; Lin, E.C.Y.; Gutowski, B., 1992a. Al_2O_3 -SiC Composites from Alumino-Silicates Precursors. *J. American Ceramic Society*, 75 (8): 2283-85.
2. Stevens, R., 1981. Zirconia: second phase particle transformation toughening of ceramics. *British Ceramic Society*, 80 (3): 81-85.
3. Becher P.F., 1986. Toughening behavior in ceramics associated with the transformation of tetragonal ZrO_2 . *Acta Metall.*, 14 (10): 1885-1891.
4. Jue, J.F. and Virkar, A.V., 1990. Fabrication, microstructural characterization, and mechanical properties of polycrystalline t'-zirconia. *J. Am. Ceram. Soc.*, 73 (12): 3650-3657.
5. Bartlett, R.W.; Wadsworth, M.E.; Cutler, I.B., 1963. The Oxidation Kinetics of Zirconium Carbide. *Transactions of the Metallurgical Society of AIME*, 227: 467-472.
6. Kuriakose, A.K. and Margrave, J.L., 1964. The Oxidation Kinetics of Zirconium Diboride and Zirconium Carbide at High Temperatures. *J. Electrochem. Society*, 111 (7): 827-831.
7. Penugonda, M.R. and Chaklader, A.C.D., 1989. Alumina-SiC Composites from Kaolinite-Carbon Precursors by Hot-Pressing. *Solid State Phenomena*, 8 (9): 457-70.
8. Chaklader, A.C.D., Lin, E.C.Y., 1993. Synthesis of Ceramic-Ceramic Composite Powders from Natural Mineral Precursors. *J. Materials Synthesis and Processing*, 1(3): 145-152.
9. Funahashi, T.; Ueda, K.; Uchimura, R.; Oguchi, Y., 1988. High-Purity Zirconia from Zircon by Carbothermic Reduction under Reduced Pressure. *Kawasaki Steel Technical Report*, 18: 73-80.
10. Shimada, S. and Ishii, Y., 1990. Oxidation Kinetics of Zirconium Carbide at Relatively Low Temperatures. *J. American Ceramic Society*, 73 (10): 2804-808.
11. Chaklader, A.C.D.; Gupta, S.D.; Lin, E.C.Y.; Gutowski, B., 1992b. Al_2O_3 -SiC Composites Using Alumino-Silicates (Natural and Synthetic) Precursor Materials. *Solid State Phenomena*, 25-26: 429-436.
12. Tamura, K.; Ogawa, T.; Fukuda, K., 1990. The Oxidation Behavior of ZrC Coating and Powder Studied by Laser Raman Spectroscopy and X-Ray Diffraction. *Journal of Nuclear Materials*, 175: 266-269.
13. Wang, J. and Stevens, R., 1988. Toughening mechanisms in duplex alumina-zirconia ceramics. *Journal of Materials Science*, 23: 804-808.
14. Uchiyama T.; Inoue, S.; Niihara, K., 1991. Multiple Toughening in Al_2O_3 /SiC Whisker/ ZrO_2 Composites. *Silicon Carbide Ceramics-1*. Ed. S. Somiya and Y. Inomata. Elsevier, London p. 265- 274.

High Temperature Flow Stress Model and Hot Deformation Behaviors of High Mo Austenitic Stainless Steel

Xu Yourong, Chen Liangshen, Jin Lei, Wang Deying

School of Materials Science&Engineering, Shanghai University,
Jiading, Shanghai, 201800, P.R.China

ABSTRACT

Single stage and double stage interrupted hot compression tests on the Thermecmaster-Z simulator for physical simulating hot rolling have been carried out for a 00Cr20Ni18Mo6Cu[N] austenitic stainless steel under high temperature range from 1223K to 1373K and various strain rates. The dynamic and static mechanical behaviors and microstructure evolution of the steel were studied. The deformation activation energies of dynamic, static and metadynamic recrystallization were calculated. A series of perfect flow stress models considering dynamic recrystallization were established. The predicted results agree well with the experiment data. Kinetic modeling of metadynamic and static recrystallization has also been determined.

Key Words: Flow stress, Hot deformation, Stainless steel

INTRODUCTION

The alloying of high molybdenum and high nitrogen is the fundamental of developing modern highly-alloyed austenitic stainless steels, its optimum resistance to pitting and stress corrosion are commonly acknowledged [1,2]. In this system, 00Cr20Ni18Mo6Cu[N] and 00Cr20Ni18Mo5Mo[N] are widely applied to oceanics, petroleum chemistry industry etc.[1,2]. However, these kinds of materials of high alloy are hard to processing and normally have worse plasticity. It is necessary to improve the processing properties during hot working and to control the microstructure.

For improving the workability and studying characterization of hot deformation of stainless steels, more considerable studies were carried out on conventional type 300 stainless steels[3-4], but systematic studies on high temperature deformation and evolution of microstructure were relatively few for high molybdenum and high nitrogen highly-alloyed austenitic stainless steels[5,6].

This paper is aimed at promoting high molybdenum and high nitrogen 00Cr20Ni18Mo6Cu[N] Austenitic stainless steel's hot workability and improving hot processing behaviors. Single and double compression tests were carried out on Thermecmaster-Z simulator, together with studying evolution of microstructure with OM, SEM and TEM, the flow stress model, hot deformation behaviors such as dynamic and static recrystallization, softening behaviors and microstructure evolution for high Mo and high N 00Cr20Ni18Mo6Cu[N] Austenitic stainless steel were investigated.

EXPERIMENTAL PROCEDURE

The chemical composition of the experimental Austenitic Stainless Steel 00Cr20Ni18Mo6Cu[N] is shown in Table 1 in weight percent.

Table 1 The Chemical Composition of 00Cr20Ni18Mo6Cu[N] steel (wt%)

	C	Si	Mn	P	S	Ni	Cr	Mo	Cu	N
wt%	<0.03	<0.80	<0.80	<0.03	<0.03	17.50-20.05	19.0-21.0	5.2-6.3	0.5-1.0	0.15-0.25

The experimental materials underwent prior homogenous solid-solute treatment. Then machined to cylinder specimens with $\varnothing 8\text{mm} \times 12\text{mm}$. The experimental specimens were compressed on the Thermecmaster-Z

simulator, using the methods of single-stage and double-stage hot compression tests. Single-stage tests were conducted under various temperature (950-1100°C) and strain rates (0.1s⁻¹ - 60s⁻¹) with strain of 0.8. Double-stage tests were conducted under various temperature (900-1050°C) and prestrain (0.3-0.6) as interpass time ranged from 1s to 200s. Hot deformed materials were rapid-cooled in Ar, N₂ or water to investigate the evolution of microstructures.

RESULTS AND DISCUSSION

Calculation of activation energy for hot deformation

Plastic deformation of materials was a process of thermo/mechanical activation. The relationship between flow stress and temperature, strain rate is showed by the following equation:

$$\dot{\epsilon} = A[\sinh(\alpha\sigma_{ss}^*)]^m \exp(-Q/RT) \quad 1.$$

$$Z = A[\sinh(\alpha\sigma_{ss}^*)]^m = \dot{\epsilon} \exp(Q/RT) \quad 2.$$

Where Q is the apparent activation energy for hot deformation, σ_{ss}^* is the saturation stress, A and m are experimental constants. α is the optimum factor which is unrelated to experimental conditions. According to McQueen, α for stainless steel is 0.012 [4]. From Equation 1 and 2, we can get the following equations:

$$Q = R[\partial \ln \dot{\epsilon} / \partial (1/T)]_{\sigma = \text{const}} \quad 3.$$

$$Q = R[\partial \ln \sinh(\alpha\sigma_{ss}^*) / \partial (1/T)]_{\dot{\epsilon}} / [\partial \ln \sinh(\alpha\sigma_{ss}^*) / \partial \ln \dot{\epsilon}]_T \quad 4.$$

The activation energy for 00Cr20Ni18Mo6Cu[N] Austenitic Stainless Steel is calculated as 586.7KJ/mol. Calculation result shows that experimental steel has much higher Q than normal low-carbon steels (200-300KJ/mol), alloyed steels and other austenitic stainless steels (400-500KJ/mol), it is due to their high alloy content (44%). The increase in alloy content, especially Mo will obviously raise the activation energy.

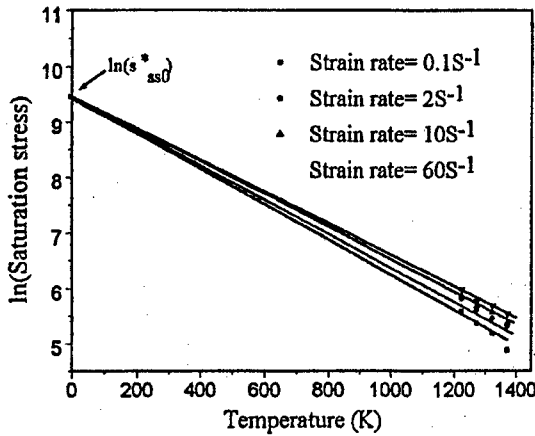


Fig. 1. Plot of σ_{ss}^* vs. T . Convergence of the extrapolation at 0 K determines the maximum saturation stress.

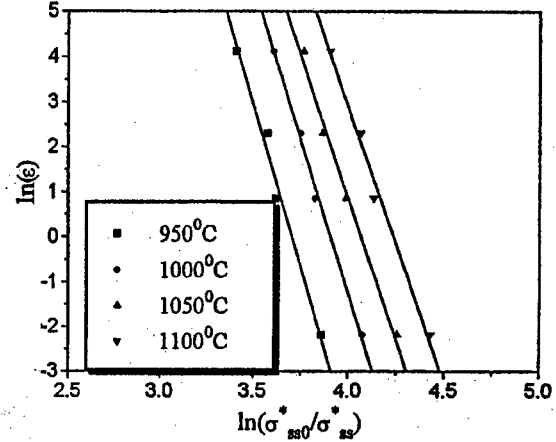


Fig. 2. Relationship between $\ln(\dot{\epsilon})$ and $\ln(\sigma_{ss0}^*/\sigma_{ss}^*)$.

ΔH can be calculated from the slopes.

Saturation Stress σ_{ss}^*

Saturation stress is the revised value of peak stress (σ_p) considering dynamic recovery (DRV) without the occurrence of dynamic recrystallization (DRX), so σ_p in the true stress-strain curves is lower than σ_{ss}^* . Dependence of σ_{ss}^* on deformation can be demonstrated by the following Kocks-Mecking equation [7]:

$$\dot{\epsilon} / \dot{\epsilon}_0 = (\sigma_{ss}^* / \sigma_{ss0}^*)^{\Gamma} \exp(-Q/RT) \quad 5.$$

$$\Delta H(\sigma) = \Gamma \ln(\sigma_{ss0}^* / \sigma_{ss}^*) \quad 6.$$

where $\dot{\epsilon}$ is strain rate, σ_{ss}^* is saturation stress, $\dot{\epsilon}_0$ is constant for the given materials, σ_{ss0}^* is the maximum saturate stress at 0 K, which can be obtained by extrapolating the $\ln(\sigma_{ss}^*)$ vs. T lines (as in Fig. 1) to 0 K. In the experiment $\sigma_{ss0}^* = 1.26 \times 10^4$ Mpa, Γ is the fault stacking energy parameter, ΔH is activation enthalpy. For a given material, ΔH varies with deformation conditions at low temperature while at high temperature, ΔH agrees reasonably well with the activation energy, ΔH can be calculated from the following equation:

$$\ln(\dot{\epsilon}) = \ln(\dot{\epsilon}_0) - \Gamma/RT \ln(\sigma_{ss0}^*/\sigma_{ss}^*) \quad 7.$$

From the $\ln(\dot{\epsilon})$ vs. $\ln(\sigma_{ss0}^*/\sigma_{ss}^*)$ plot shown in Fig. 2., Γ is determined. According to Equation 6, ΔH can be obtained easily. In the experiment, ΔH increases as temperature increases and decreases as strain rate increases -- the values range from 480 to 620, and approximate to the activation energy at high temperature.

Kinetics of dynamic recrystallization

The fraction of dynamic recrystallization is defined as:

$$X = (\sigma_{ss}^* - \sigma) / (\sigma_{ss}^* - \sigma_s) \quad 8.$$

Where σ_{ss}^* is the saturation stress, σ_s is the steady flow stress, σ is the flow stress under various deformation conditions. According to the Johnson-Mehl-Avrami equation:

$$X = 1 - \exp[-A(\epsilon - \epsilon_p)^n] \quad 9.$$

Where n is the Avrami coefficient, A is a constant, ϵ_p is the peak strain, former investigation shows [8]:

$$\epsilon_p = K' Z^m = 1.84 \times 10^{-4} Z^{0.128} \quad 11.$$

m and K' are calculated as $m = 0.128$, $K' = 1.84 \times 10^{-4}$. The linear relationship of $\ln(\ln(1/(1-X)))$ vs. $\ln(\epsilon - \epsilon_p)$ is shown in Fig. 3, but the lines are not parallel to each other, demonstrating that the Avrami coefficient (n) is not a constant under various deformation conditions. Under experimental conditions, the value of n ranges from 0.9 to 2, depending on deformation parameters. The mathematical model for the Avrami coefficient is established as follows:

$$n = A \dot{\epsilon}^p \exp(Q/RT) = 0.048 \dot{\epsilon}^{0.061} \exp(35.99(\text{KJ/mol})/RT) \quad 12.$$

A , p , Q are calculated as $A = 0.048$, $p = 0.061$, $Q = 35.99 \text{ KJ/mol}$ respectively.

There are different views concerning the Avrami coefficient n . Larsraoin [9] held that the value of n is a constant of about 1.5, but this idea does not conform with reality. Febregue [10] considered n to be a function of $\dot{\epsilon}$ and T . The results of our experiments are similar to that of Febregue. There is great variation in the value of coefficient A . In the experiment, the value of A ranged from 1.49 to 8.17. Generally, it is considered that A depends on deformation conditions, but this needs further research.

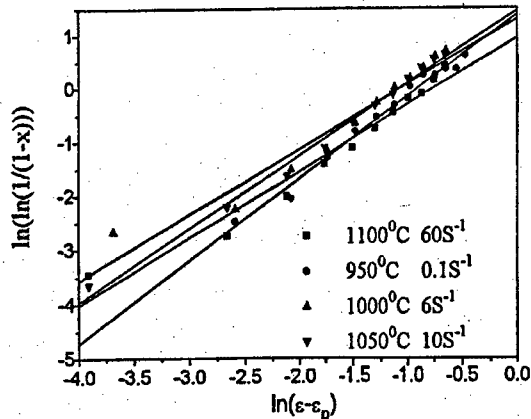


Fig. 3. Relationship of $\ln(\ln(1/(1-X)))$ vs. $\ln(\epsilon - \epsilon_p)$

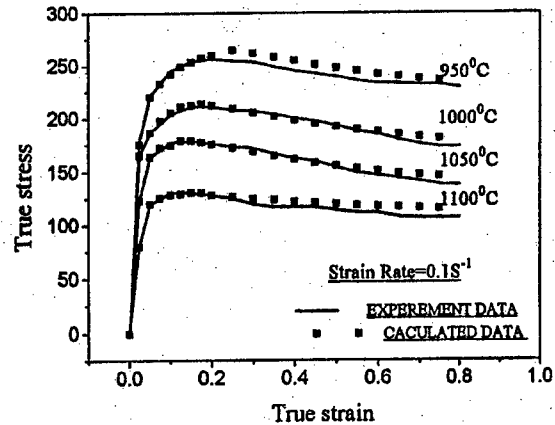


Fig. 4. Comparison of model prediction and experimental data.

The Flow Stress Model Before and After Dynamic Recrystallization

Strain-Stress curves show the internal link of flow stress and deformation conditions and the evolution of microstructures. Here using the Jonas [11] model to define the flow stress before peak stress:

$$\sigma^2 = \sigma_{ss}^{*2} + (\sigma_0^2 - \sigma_{ss}^{*2})e^{-\Omega\epsilon} \quad 13.$$

where σ_0 is the initial stress and Ω is the parameter of dynamic recovery. The initial stress, σ_0 , can be obtained by means of regression or by reading directly from strain-stress curves. The value of Ω depends on the deformation temperature and strain rate. For the experiment, Ω was defined as:

$$\Omega = 4.42 \times 10^5 \dot{\epsilon}^{-0.093} \exp(1.08 \times 10^5 / RT)$$

which agrees with the result of Yoshie's experiment [12].

The relationship of σ_{ss}^* and deformation conditions has been investigated in this paper. From Equation 1 and 2, σ_{ss}^* can also be expressed as:

$$\sigma_{ss}^* = \sin^{-1} h(AZ^{1/m}) / \alpha 83.3 \sin^{-1} h(2.17 \times 10^{-5} Z^{0.237}) \quad 14.$$

where A and 1/m are calculated as $A = 2.17 \times 10^{-5}$ and $1/m = 0.237$, respectively.

A softening function is not considered in the model before dynamic recrystallization, while after dynamic recrystallization, it must be taken into account. By combining Equations 8 and 9, σ can be described as:

$$\sigma = \sigma_{ss}^* - (\sigma_{ss}^* - \sigma_s) \{1 - \exp[-A(\epsilon - \epsilon_p)^n]\} \quad 15.$$

So, the integrated mathematical flow stress models have now been formulated. To verify reliability of the model, we compared the predicted data with the experiment data in Fig. 4., which demonstrates the data are approximately in accord with each other, so the model reliability is confirmed.

Evolution of Microstructures for Dynamic Softening

Because of the high alloy ingredients in the experimental steel, dynamic recrystallization is retarded, therefore, the steel has a high recrystallization temperature. Fig. 5 represents the effects of deformation temperature and strain rate on dynamic recrystallization. If the temperature is only 950 °C, even with a strain rate as high as $60s^{-1}$ (high strain rate can hasten recrystallization), only a small quantity of recrystallized grains appear at the grain boundaries. When the temperature is below 950 °C, it is difficult for recrystallization to occur. An increase in temperature can promote recrystallization. When the temperature reaches 1000 °C, a large quantity of recrystallized grains appear and expand into the initial crystal as the strain rate is only $0.1s^{-1}$. The dynamic recrystallization tends to complete at high strain rate as temperature increases up to 1050 °C, it is indicated that high strain rate can enhance recrystallization at high temperature.

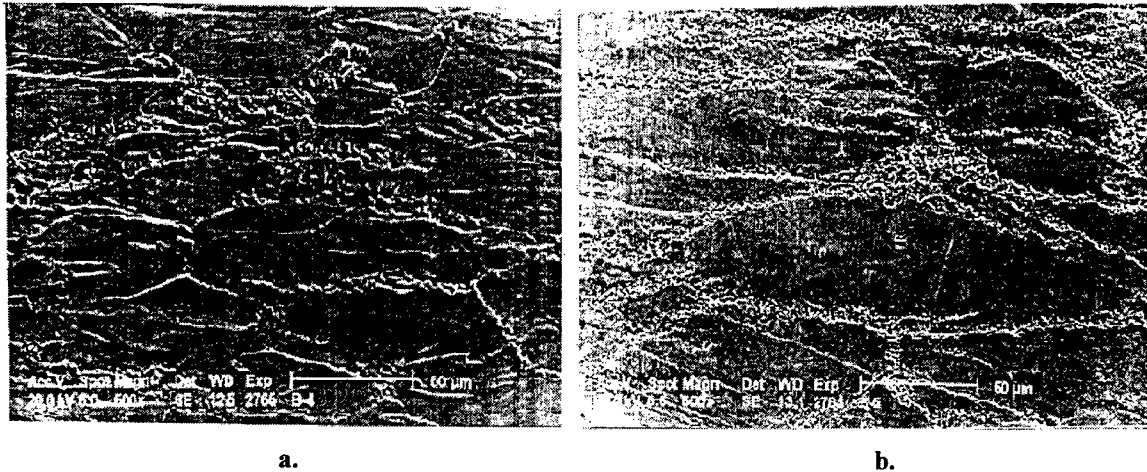


Fig. 5. Microstructures of dynamic modelling (SEM) a. $T=950\text{ }^{\circ}\text{C}$, $\dot{\epsilon}=60s^{-1}$ b. $T=1000\text{ }^{\circ}\text{C}$, $\dot{\epsilon}=0.1s^{-1}$

Metadynamic Recrystallization and Static Recrystallization Results

By interrupting the compression tests beyond the critical strain of dynamic recrystallization, metadynamic recrystallization (MDRX) and static recrystallization (SRX) takes place during the interpass time. 75% softening occurred after 100 seconds of maintaining the prestrain at 0.4, the strain rate at 2 s^{-1} , the holding temperature at 1050°C . Elongating the interpass time from 100s to 200s caused the softening fraction to increase only slightly as shown in Fig. 6. It is demonstrated that the softening process is complete after 100 seconds of interpass, but work-hardening was not completely extinct. Incomplete softening deduced from the appearance of MDRX was reported by Sakai et.al. [13] from investigations into hot deformation behavior of Ni multi-crystals.

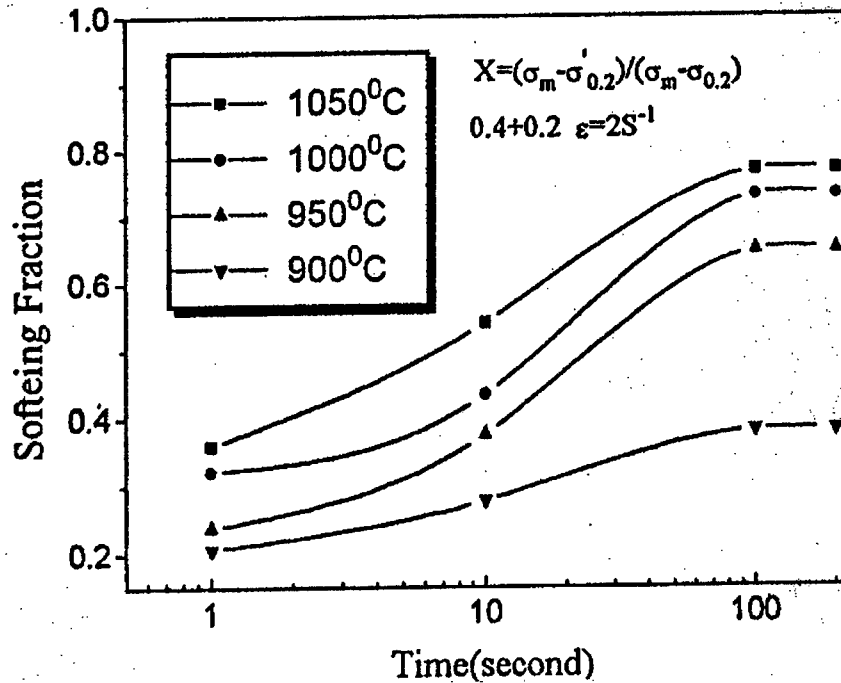


Fig. 6. Static softening fraction curves of high Mo austenitic stainless steel

The kinetics of conventional static recrystallization can be described by the following equation [14]:

$$X = 1 - \exp[-0.693(t/t_{0.5})^n] \quad 16.$$

where X is the recrystallized volume fraction, n is the Avrami exponent, (about 0.5 for the experimental steel), $t_{0.5}$ is the time for 50% recrystallization, which is an important parameter for SRX and MDRX which implies the rate of SRX and MDRX. As for SRX, prestrain greatly influences the value of $t_{0.5}$, while strain rate shows little effect on $t_{0.5}$. In fact, most researchers simply choose to ignore the effect of strain rate on SRX. But for MDRX, the value of $t_{0.5}$ largely depends on strain rate, while the effect of prestrain can be ignored [15]. So the following equations are obtained for SRX and MDRX, respectively:

$$t_{0.5} = A_1 \epsilon^m \exp(Q_{\text{SRX}}/RT) \quad , \quad (\text{For SRX}) \quad 17.$$

$$t_{0.5} = A_2 \dot{\epsilon}^n \exp(Q_{\text{MDRX}}/RT) \quad , \quad (\text{For MDRX}) \quad 18.$$

where Q_{SRX} , Q_{MDRX} are activation energies for SRX and MDRX, respectively. By nonlinear regression, the value of the experimental constants for high-Mo stainless steel was determined to be:

$$A_1 = 2.21 \times 10^{-11}, \quad A_2 = 5.60 \times 10^{-9}, \quad m = -1.79, \quad n = -0.365, \quad Q_{\text{SRX}} = 483.7 \text{ KJ/mol}, \quad Q_{\text{MDRX}} = 253.5 \text{ KJ/mol}.$$

Because metadynamic recrystallization is a growth process of a dynamic recrystallized nucleus, it has no pregnant period, so the activation energy is much lower than that of static recrystallization.

CONCLUSION

1 The activation energies for hot deformation, static recrystallization and metadynamic recrystallization are 586.7 KJ/mol, 483.7 KJ/mol and 253.5 KJ/mol, respectively.

2 The flow stress models are established as:

$$\sigma^2 = \sigma_{ss}^{*2} + (\sigma_0^2 - \sigma_{ss}^{*2}) e^{-\Omega \epsilon}, \quad (\epsilon < \epsilon_p), \quad \sigma = \sigma_{ss}^* - (\sigma_{ss}^* - \sigma_{ds}) \{1 - \exp[-A(\epsilon - \epsilon_p)^n]\}, \quad (\epsilon \geq \epsilon_p)$$

$$\sigma_{ss}^* = 83.3 \sin^{-1} h(2.17 \times 10^{-5} Z^{0.237}), \quad \Omega = 4.42 \times 10^5 \dot{\epsilon}^{-0.093} \exp(1.08 \times 10^5 / RT),$$

$$n = 0.048 \dot{\epsilon}^{0.061} \exp(3.60 \times 10^4 / RT), \quad \epsilon_p = 1.84 \times 10^{-4} Z^{0.128}$$

3 The Avrami coefficient of MDRX is about 0.5, and the kinetics of MDRX and SRX are demonstrated by following equations:

$$t_{0.5} = 2.29 \times 10^{-11} \dot{\epsilon}^{-1.79} \exp(4.83 \times 10^5 / RT), \quad (\text{For SRX})$$

$$t_{0.5} = 5.6 \times 10^{-9} \dot{\epsilon}^{-0.365} \exp(253.5 \times 10^5 / RT), \quad (\text{For MDRX})$$

4 For dynamic recrystallization to be retarded by high alloy ingredients, the tested steel must have a high recrystallization temperature. So, carrying out single-stage deformation above 1050 °C and using high strain rate (10 s⁻¹ - 60 s⁻¹) is advantageous for inducing dynamic recrystallization and decreasing cracks.

5 Increasing holding temperature and interpass time can assist in the occurrence of static recrystallization.

REFERENCES

1. R.F.A. Jargelius-Pettersson, 1996. Scandinavian Journal of Metallurgy, 24(5/6), 188.
2. Staffan, Hertzmen, 1996. Scandinavian Journal of Metallurgy, 24(4), 140.
3. N. Ryan and H.J. McQueen, 1994. "Advances in hot deformation textures and microstructures", ed. J.J. Jonas et al., The Minerals, Metals & Materials Society, 445.
4. N.D. Ryan and H.J. McQueen, 1990. Materials Forum, 14, 283.
5. A.A. Holmvik, J.K. Salberg, J. Perttula, 1997. "Thermec97", International Conference on Thermomechanical Processing of Steels and Other Materials, ed. T. Chandra and T. Sakai, 1, 241.
6. W. Roberts, 1994. "Deformation, Processing and Structure", ed. G. Krauss, ASM, 109.
7. U.F. Kocks, H. Mecking, 1985. Strength of Metals and Alloys, ed. H.J. McQueen et al., (Oxford, England: Pergamon Press, 345.
8. C. Roucoules, S. Yue and J.J. Jonas, 1995. "Microalloying95", Conference Proceeding, 165.
9. B. Dutta and C.M. Sellars, 1987. Materials Science and Technology, March, 3, 197.
10. P. Fabregue, 1994. Advances in Hot Deformation Textures and Microstructures, ed. J.J. Jonas et al., The Materials, Metals & Materials Society, 75.
11. A. Lassraoui, J.J. Jonas, 1991. Metallurgical Transactions 22A, July, 1547.
12. A. Yoshie, H. Morikawa et al., 1987. Trans. ISU, 27, 29.
13. M. Militier et al., 1994. Acta Metall Mater, 42, 133.
14. P.D. Hodgson, D.C. Collinson, B.A. Parker, 1994. Advances in Hot Deformation Textures and Microstructures, ed. by J.J. Jonas, T.R. Bieler et al, The Materials, Metals & Materials Society, 41.
15. C. Roucoules, S. Yue, J.J. Jonas, 1993. 1st Inter. Conf. on Modelling Metal Rolling Process, London, The Institute of Materials, 165.

Intelligent Manufacturing III

MANAGEMENT OF INFORMATION IN COMPLEX SYSTEMS: PERSPECTIVES FOR THE NEW MILLENNIUM

E. Szczerbicki* and Z. Gomolka**

*The University of Newcastle, Newcastle, Australia

**The University of Szczecin, Szczecin, Poland

ABSTRACT

Engineering, operations research, and management science use scientific and engineering processes to design, plan, and schedule increasingly more complex industrial systems in order to enhance performance. One can argue that the systems have grown in complexity over the years mainly due to increased strive for resource optimization combined with a greater degree of uncertainty in the system's environment. Information is seen as one of the main resources that managers try to use in an optimal way. Managing complex systems requires a greater understanding and knowledge about the role of information in systems operation. Today, a *growing complexity of information flow* is a characteristics of enterprises which concerns products to be manufactured, services to be offered, processes and company structures. Complex systems also operate in *changing environments* surrounded by numerous *uncertainties* and *disturbances*. Difficulties arise from unexpected tasks and events and from a multitude of possible failures and other interactions during the attempt to control various activities in dynamic environments. Therefore, management of information is one of the most important aspects to be considered in intelligent management systems, which are expected to solve unforeseen problems, even on the basis of incomplete and imprecise information. The paper discusses the importance of information in operation management as well as new challenges in information modelling, visualisation and communication in information society.

INTRODUCTION

Information becomes an increasingly more important resource in all kinds of business, industrial, and service operations. Changes, uncertainty, imprecision, and complexity became the most important factors affecting the behaviour of modern markets. Functioning in such markets requires increasing amount of information to be processed in substantially shorter periods of time. Therefore the time span left for decision making is dramatically decreased.

Some major problems associated with these facts can be traced to the level of operational management of a company. The efficiency of management at that level depends mainly on the amount of time needed to react to changes in both internal and external environments in which a given company is functioning. This efficiency depends heavily on the realization cycle of clients' tasks, tasks being the *spiritus movens* of any business, and the lack of them resulting in operational termination and bankruptcy.

The operational management level is increasingly more often a decisive factor in a company's survival and expansion. We have to concentrate on this level of management as much as on tactical and strategic levels that are usually the main focus of companies today. The operational level at which tasks are processed is no longer a stable one. A number of parameters associated with these tasks can change and will be changing more and more often as the dynamics of external environment increases.

Demand for Information

Functioning in uncertain and imprecise conditions requires predictions of future states of environment in which systems operate. It requires increasingly more efficient and intelligent decision support tools that are able to cope with unexpected changes. Application of such tools usually means significantly greater demand for information and larger amount of information must be processed at all management levels.

Information that is needed often originates at different, geographically distributed sources and is available in different forms and different coding. Thus, new tools are needed to cope with this emerging problem of information diversity. The challenge of the next millenium will be to retrieve and transform huge amounts of different forms of information into knowledge needed to support our decision making processes.

CHALLENGES OF THE NEXT MILLENNIUM

Managing companies in the next century, the century of information society, will necessitate the use of new means of communication with external environment. It will also require much greater adaptability of companies, it will require the companies of the next Millennium to be transformed into intelligent, learning organizations able to cope with globalization of information resources. This globalization means that the main problem will not be the access to information but the ability to mine it and then to transform it into a useful operational and strategic resource.

The increasing frequency of change in the state of the environment in which a company operates, creates an important challenge related to time. Time becomes a decisive factor in information retrieval and decision-making processes. Managing complex industrial systems (manufacturing, processing, distribution, servicing, mining, etc.) that function in uncertain information-rich environments requires greater understanding and knowledge of the role of information in systems operation. To gain this understanding, a theory will be needed that could be used to model and evaluate information flow in different situations.

In fact, our needs for the next century go well beyond the above in requirement for a theory considering important practical issues of information, i.e. delays, incompleteness, imprecision and loss in value. The current practice of dealing with such issues are mostly when problems are detected and reactively. This situation may not be desirable and definitely be a major drawback for complex systems that more and more rely on the timeliness and quality of information for their operation. A theory, in this respect, would greatly enhance the understanding of the various factors that influence the quality of information to the benefit of better decisions in adequate time.

Systems become increasingly complex. Their decomposition into smaller units is the usual way to overcome the problem of complexity. This has historically led to the development of atomized structures consisting of a limited number of *autonomous subsystems* that decide about their own information input and output requirements, i.e. can be characterized by what is called an *information closure*. Autonomous subsystems can still be interrelated and embedded in larger systems, as autonomy and independence are not equivalent concepts. These ideas are recently gaining very strong interest in both academia and industry, and the atomized approach to information flow modelling and evaluation is an idea whose time has certainly come. [1, 2, 3, 4].

In a real-world context autonomous subsystems consist of groups of people and/or machines tied by the flow of information both within a given subsystem and between this subsystem and its external environment [5, 6]. A theory is needed that could be used to evaluate such an information flow. The theory should allow for the evaluation of an information flow to be performed for different types of external and internal environments of a given subsystem. It should take into account two basic cases, i.e., static and dynamic processes describing the external environment. Such issues as the role of correlation and interaction, and the losses caused by incomplete and delayed information should also be considered. The theory should also accommodate the question of uncertain and imprecise information flow modelling.

In particular we will have to address the frequent situations in which the following should be answered:

"How to structure an exchange of information between a system and its uncertain, dynamic and imprecise environment?"

and

"What is better, complete information but heavily delayed, or incomplete information less delayed?"

The value of information that flows within a given subsystem is different for different information structures and different environments [6, 7]. It can be considerably affected by two major attributes of information: incompleteness and delay. The highest value will be possessed by a full information structure (including all relevant information possible). On the other hand, gathering information in a dynamic environment causes its delay. Both delay and incompleteness can be represented by losses in the value of the information structure. Currently, there is no theoretical foundation for such a representation but managers of the next millennium will certainly need it.

The delay of information combined with the dynamics in the environment can cause substantial losses in the value of information as a useful resource in decision support. We have to turn huge amounts of information into knowledge needed for our knowledge-based systems very fast. Quick perception of information becomes an important issue. Another challenge emerges here - visualization of information.

Visualization of Information

We no longer have time to study pages of reports and columns of data. We have to visualize information quickly and effectively. New tools are needed to support the ways we communicate information to the decision maker. Visualization and presentation of information becomes one of the most important areas of research in Cybernetics and Artificial Intelligence.

Visualization can help make sense of the flood of data. When applied with some insight into visual perception and with attention to the nature of the data, and how the data are to be used visualization can become a very powerful tool in future intelligent information systems. Current and future research trends in this area include such important topics as [8, 9, 10, 11, 12]:

- colour and information,
- complexity and clarity of human perception,
- use of multi-media, internet, and WWW screens,
- animation for scientific visualizations,
- design of efficient computer interfaces.

Use of colour in presenting data becomes an increasingly important research topic. The key issue is the color-map which may be defined as a mapping from data to colour. In most colour maps, red is mapped to the highest data value, blue to the lowest, and the other data values are interpolated along the full extent of the rainbow spectrum. An example would be a temperature profile over land mass on a weather map.

But there are some unsolved problems related to colour data representation. Colour is a perceptual phenomenon. What is commonly called colour is only one of three parameters. Another is the brightness of the signal - intensity. The third is the admixture of white - saturation. To add to the complexity of the problem, the above parameters' relationship to what is perceived is nonlinear. Colour perception issue in data mapping is one of the challenges in information management for the next century.

Opportunities and Challenges

Clearly, better visual representations of data are needed, particularly by way of colormaps that will induce more faithful impressions of the structure hidden in the data. The basic challenge here is: how should colour be used to encode characteristics of interest in a dataset. Perceptual encoding becomes an area of focus in visualization of information. Some pioneering work here has already been done at IBM for its product known as Visualization Data Explorer (DX) which is a visual language, object based package [8]. A software called Pravda (perceptual rule-based architecture for visualizing data accurately) is a support tool developed to choose for colormapping based on principles of perception and colour theory. Pravda can be used interactively with DX. The data to be visualized are imported into DX and flow in to a module called PravdaColor. This tool determines the data's characteristics including their spatial frequency. With the aid of a control panel, the user can select the colormapping goal of the final visualization. Pravda is not available commercially, but a version of it may be in the future.

The ability of presenting information in a picture form becomes crucial as this form is the fastest and most natural way of communicating data to a decision maker. This form combined with voice, animation and colour and presented by multimedia techniques as *ars electronica* transforms real-life reality in which a company functions into virtual reality in cyberspace. Management of information in the next millenium will focus a good part of its efforts on a number of cyberspace related techniques and tools and their implementation in decision support processes.

Of some interest here may be the fact that one of the best graphical representations of information was developed in 1869 by Charles Minard [12]. It depicts the losses suffered by Napoleon's army in the Russian campaign of 1812. Six variables were plotted in that graph which tells a rich coherent story with its

multivariate data, far more enlightening than just a single number posed again time. Minard's chart is still regarded as one of the best graphical representation of data ever drawn.

Presentation of information in cyberspace, the virtual reality perspective from which we will be able to observe and judge how a given complex system functions, creates new opportunities for better, more efficient management of companies of the next century. There also are, however, completely new dangers that we have to be aware of.

The time factor is one of the main decision making constraints in the sense that we cannot go back and change decisions that were not optimal. Virtual reality often develops an illusion that the above is possible. This may cause serious underestimation of risks associated with our decisions. Another problem associated with multimedia forms of information visualization is its security. New tools will be needed to restrict access to confidential parts of information.

As computer systems processing information in increasingly more dynamic and uncertain environments become more powerful and complex, our interactions with them have become more information laden and, consequently, more burdensome. It is now generally recognized that new intelligent user interfaces will be needed. The pieces of solutions to this problem are coming together from a variety of disciplines, including machine learning, user modeling, intelligent tutoring, information retrieval, and data mining. Furthermore, related work is discussed in the field of autonomous multiagent systems.

Despite the above dangers the interest in visualization of information and its presentation through various forms of *ars electronica* is growing very rapidly. It is seen as one of the most important means of possible increase of efficiency and quality of information management in the next millenium.

REFERENCES

1. A. Gunasekaran, M. Sarhadi, 1997. Planning and management issues in enterprise integration, Concurrent Engineering: Research and Application.
2. L. Pacholski, M. Wejman, 1995. Soft Modelling of the Ergonomics of the Multiagent Manufacturing Systems, Taylor and Francis.
3. J. Raczekowski, W. Reithofer, 1998. Design of Consistent Enterprise Models, Cybernetics and Systems: An International Journal.
4. A. Tharumarajah, 1998. A self-organising model for scheduling distributed autonomous manufacturing agents, Cybernetics and Systems: An International Journal.
5. Z. Gomolka, 1995. Elements of General Systems Theory and Systems Modelling, WNUS.
6. E. Szczerbicki, 1992. Information Processing for the Development of Integrated Multiagent Manufacturing Systems, Gdansk University Press.
7. Z. Gomolka and E. Szczerbicki, 1997. Cybernetics of a goal-seeking agents: conceptual model, Systems Analysis, Modelling, Simulation.
8. B.E. Rogowitz and L. Treinish, 1998. Data visualization: the end of the rainbow, IEEE Spectrum.
9. H. Lefkowitz and G.T. Herman, 1992. Color scales for image data, IEEE Computer Graphics and Applications.
10. E. Tufte, 1990. The visual display of quantitative information, Graphics Press.
11. E. Tufte, 1994. Envisioning information, Graphics Press.
12. E. Tufte, 1998. Visual explanations, Graphics Press.

Present Status of Intelligent Machines in Sheet Metal Fabricating and Forming in Japan

Jun-ichi Endou

Kanagawa Institute of Technology, Atsugi, Kanagawa, Japan

Email: endo@ctas1.mse.kanagawa-it.ac.jp

ABSTRACT

The need for intelligent machines in sheet metal fabricating is explained and the concept of an intelligent machine is introduced. Some examples of intelligent bending machines are introduced and their drawbacks discussed. Intelligent machines for sheet metal forming are also introduced.

INTRODUCTION

The manufacturing circumstances of Japanese industries, especially the sheet metal fabricating industries, is very difficult these days. Small enterprises in this field are faced with problems such as high labor costs, diversification of customers needs and an aging labor force. Skills are not easy to transfer to new young laborers because of the difficulties in employing young workers these days in this industry. Some companies face the crisis of continuing their business because of these problems.

Intelligent manufacturing methodologies are expected to deal with this critical circumstance. Fully-automated intelligent machines which produce products autonomously are a dream and the final objective of engineers and scientists in the field of manufacturing science. If these enterprises can use intelligent sheet metal fabricating machines, they will be able to continue their functions profitably. The methodology, however, is not always applicable and useful in practical problems. This paper presents the present status of intelligent machines in sheet metal fabricating and forming in Japan.

STATUS OF INTELLIGENT MACHINE IN SHEET METAL FABRICATING

The Concept of the Intelligent Machine

The concept of the intelligent machine has not yet been established in a standardized way. The author has proposed a concept of the intelligent machine analogous with the human being [1] (see Figure 1.).

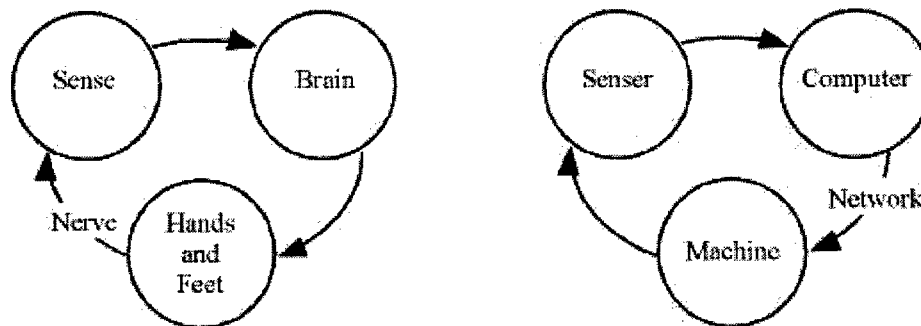


Fig. 1. Notion of the intelligent machine.

An intelligent machine must have a computer instead of a head and brain; sensing devices instead of senses, i.e., eyes, nose, etc.; mechanical working functions instead of hands and feet; and a network instead of a central nervous system which connects these functions organically. Special attention should be paid to the sensing devices that are not only for feedback control such as a tachometer of a servo-motor attached for closed loop NC, but also for higher optimization of operations.

Intelligent Bending Machines

Sheet metal fabricating consists of sheet metal shearing/cutting and punching, bending and joining / welding / connecting. Tapping and deburring are done during these processes, but these are outside the scope of the present paper. Some bending machines are known as intelligent machines. Bending processes, especially bending by press-brake, is one of the most flexible processing operations and so, bending machines require operators with special skills. Figure 2 shows a typical press-brake. Processing times of bending by press-brake are usually long. These characteristics of the press-brake has lead to the development of it as an intelligent machine.

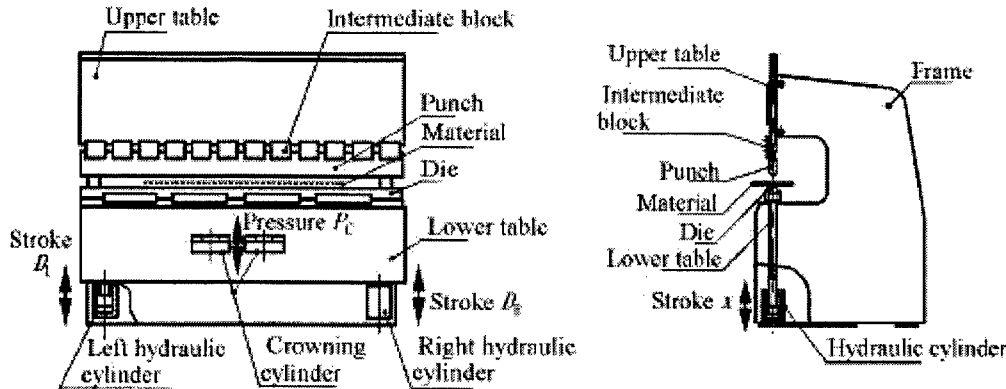


Fig. 2. View of up-stroking press-brake.

Investigations and development of intelligent bending machines in universities [2,3,4] have been conducted over the past ten years in Japan. The key feature of these developed machines is to decide operating conditions by themselves from measuring and identifying the material characteristics. Figure 3 shows an example of an intelligent bending system developed by Yang et. al.[4]. Bending load, stroke and bending angle of the metal sheet are measured and the material characteristics are identified. The correct operating conditions are then decided by the computer to determine the objective bending angle.

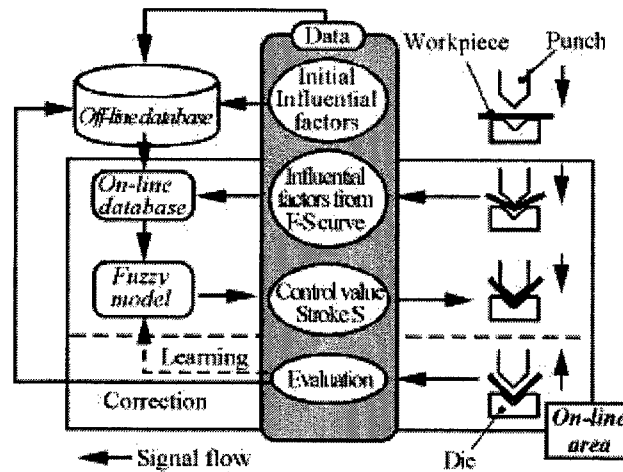


Fig. 3. Block chart of intelligent bending system.

The methods used to identify material characteristics are as follows:

- 1) simulation or theoretical analysis [2].
- 2) utilization of a data base [3,4].

If we consider the calculation time of a computer, the latter method is considered to be more effective. However, it is usually difficult to predict bending deformation of a specimen perfectly by utilizing a data base alone. In the system in Figure 3, Fuzzy Control is adopted to correct the deformation predicted by the data base. This provides a method to rapidly approximate the correct deformation requirements. To construct a data base for bending deformation of many materials requires many bending tests. Simulation

or theoretical analysis is sometimes useful to help to construct the data base. In simulation or theoretical analysis, material characteristics are easily changed, but we must note that simulated or calculated behavior is sometimes isolated from the actual bending behavior.

Developments of intelligent press-brakes in industry has also been done. A primitive intelligent press-brake was put on the market over ten years ago. The machine is driven by an AC-servo motor and has in-process sensing functions for load and stroke. By utilizing a load-stroke diagram, this machine has good reproducibility compared with conventional press-brakes. Using almost the same idea, hydraulic press-brake have also been developed with in-process load-sensing devices [5]. In Figure 4a, the typical relationship between bending forces, stroke and time are shown, while in Figure 4b, schematic diagrams of the deformation of bending, i.e., air-bending and bottoming, are shown. The maximum value of loads in the air-bending domain (F_m in Figure 4a) has positive correlation with the thickness of the metal sheet. The algorithm of this newly-developed intelligent press brake is as follows: Reference maximum bending loads in the air-bending domain is measured and stored in the computer. A slight change in the thickness of the metal sheet is identified by the slight deviation in the maximum load, and by utilizing this deviation, the final bending load in the bottoming domain is calculated.

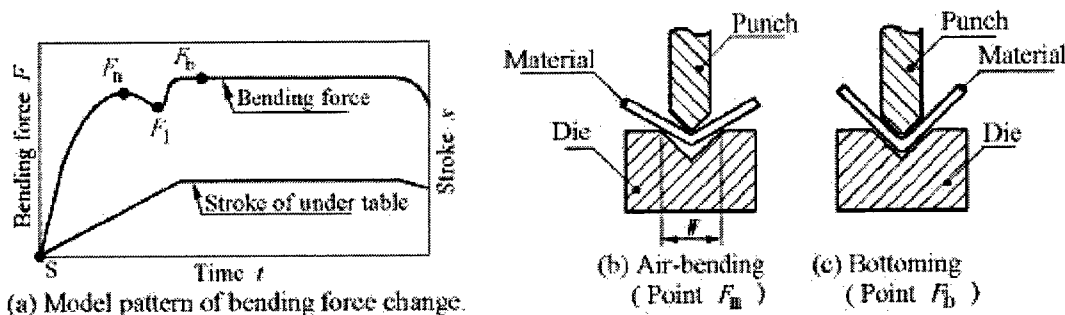


Fig. 4. Model pattern of bending force change (a), and diagrams of the deformation of bending (b & c).

Figure 5 shows a comparison of scattered bending angle of conventional press brakes (white squares in Figure 5) with the results of newly-developed press brake (black circles in Figure 5). Improvement in reproducibility of bending angle can be clearly seen as compared with the original machine. A method which improves the accuracy of longitudinal bending angle has also been developed [6]. In the bending of long metal sheet, maintaining the accuracy of the bending angle in the longitudinal direction is difficult for ordinary press-brakes because of their structure. A so-called "crowning" method which tries to control the profile of the upper and lower tables of the press-brake is adopted in order to improve accuracy of the longitudinal bending angle. The new algorithm is as follows: Pressure sensors are buried in the intermediate blocks (see Figure 2), and the pressure distribution during bending is controlled to maintain uniform controlling pressure of the crowning cylinders and hydraulic cylinders (see Figure 2).

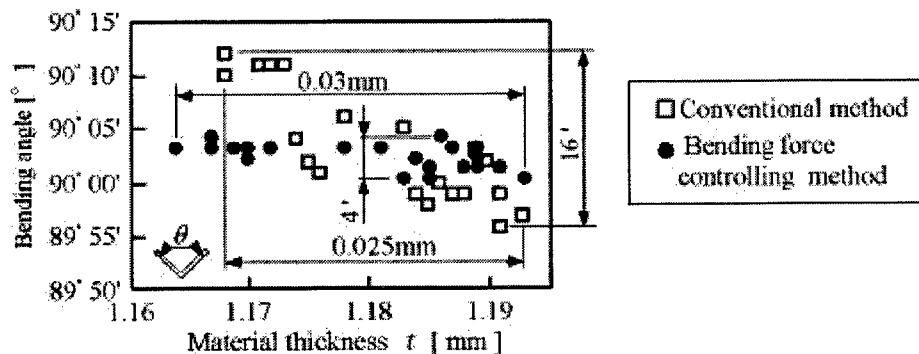


Fig. 5. Comparison of bending angle by conventional method with by newly developed method.

Figure 6a and 6b show the longitudinal distribution of bending angle and bending force ratio measured by pressure sensors with the conventional method, (Figure 6a), and with the new algorithm (Figure 6b). Considerable improvement can be seen.

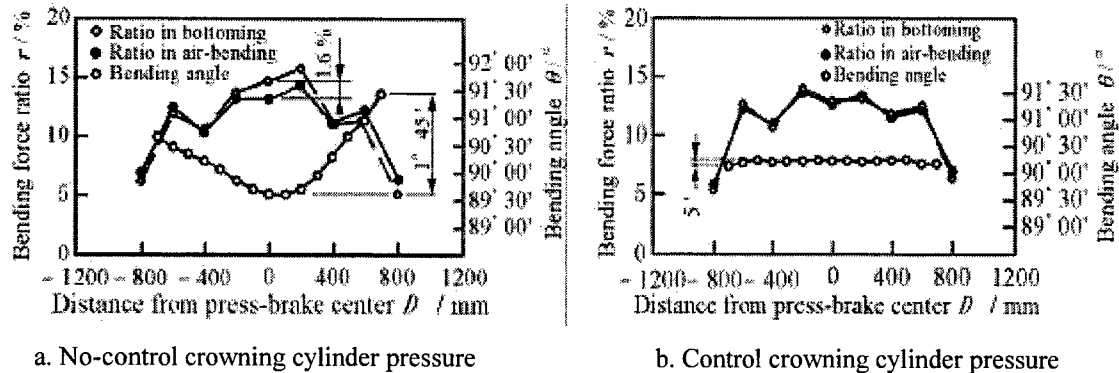


Fig. 6. Bending force ratio in bottoming and air-bending, and bending angle ($L=1600$ mm).

In-process sensing devices for bending angle measurement are also important. An image-processing-type bending-angle sensor has been developed [7]. This sensor uses a laser with the slit beam projected onto the work-piece surface. The image of the projected line on the work-piece is recorded by a CCD camera. Pattern recognition is used to calculate the bending angle with extreme accuracy. A press-brake machine with this bending angle sensor is already available on the market.

Drawbacks of Newly Developed Intelligent Press-Brake and Angle-Sensing Devices

These newly developed intelligent machines do have drawbacks. The accuracy of press-brake bending can be classified as follows:

1. Accuracy of the bending angle (Deviation between bent angle and objective angle).
2. Reproducibility of the bending angle.
3. Accuracy of the longitudinal bending angle (Longitudinal distribution of bent angle).

The purpose of the development of the first example mentioned above by Yang et al., was to improve the accuracy of the bending angle by using intelligent methodologies, i.e., to decide operating conditions for the objective angle autonomously for any material and to decrease the deviation between the bent angle and the objective angle. In order to attain this purpose, the developed system must have many sensing devices and hardware for deciding material characteristics and so on. The more sensors, the more expensive is the machine. If the cost is too high the developed intelligent bending machine can not be put on the market effectively.

On the other hand, the second and third examples are aimed at the reproducibility and accuracy of longitudinal bending angle respectively by using intelligent methodologies. These new bending machines have no problem with cost, but because of cost-reduction, their functions are limited such that a lack of functions causes drawbacks. For example, the hydraulic press-brake cannot bend metal sheet accurately when the material of the metal sheet does not have its maximum bending load in the air-bending domain. Materials which follow a linear hardening law do not have their maximum bending load in the air-bending domain. Some types of stainless steel have such characteristics. In order to avoid this drawback, a sensing device for measuring the bending angle must be attached increasing the cost of the machine.

Bending angle sensors are not always useful, necessary or accurate. Whether the sensor is a contact or non-contact type, it is not able to measure bending angle accurately when the work-piece has protrusions. Mechanical parts produced by sheet metal fabricating, sometimes have protrusions, holes, etc., and these cause unreliable bending angle measurements. Therefore, it is important that the function and cost of intelligent machines are balanced.

INTELLIGENT MACHINES FOR SHEET METAL FORMING.

Manabe et al. developed an intelligent deep-drawing machine that can optimize the blank holder pressure [8]. Figure 7 shows the developed system. In the early stages of deep drawing, material characteristics are identified from the measured values of punch force, punch stroke and blank movement, in order to optimize

the blank holder pressure. It is rather easy to decide optimum blank holder pressure for deep drawing of a cylindrical cup. But for an arbitrary figure, the problem relate to determining the optimum blank holder pressure or even, whether an optimum blank holder pressure exists.

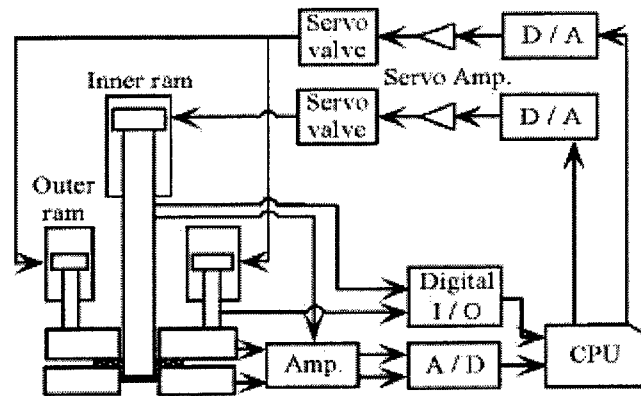


Fig. 7. Block chart of intelligent deep drawing system.

Another drawback of the intelligent deep-drawing system is its speed. A hydraulic press is used, but the speed is lower than in a conventional mechanical press. This means that productivity of the system is low and will not be accepted by the market.

Kawai et al. published on a newly-developed spinning machine with adaptive control [9]. Figure 8 shows the spinning system. In this machine, the spinning force and surface condition are sensed and operating conditions, i.e., tool path, spinning force, feed speed and spindle speed, are optimized by a data-base. The data base was well-constructed since the operating data were gathered by a research committee of JSTP.

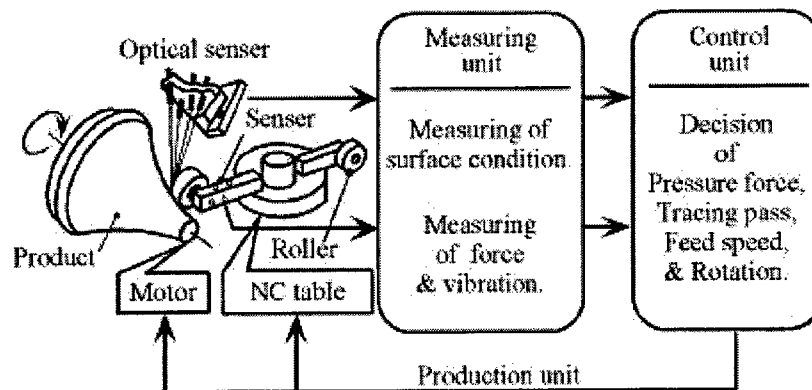


Fig. 8. Block chart of intelligent spinning system.

CONCLUSIONS

The need and concept of intelligent machines in metal forming have been described and the present status of intelligent machines for sheet metal fabricating and forming and their drawbacks have been explained.

One important question is as follows: " Will an intelligent machine be able to replace an expert operator ? "

The author doesn't think so, because the computer does not have creative faculties while expert operators do. Man and intelligent machine should be complemented to fit with each other.

REFERENCES

1. J. Endou, 1990, Expectations and problems in intelligent methodology in FA (in Japanese), J. JSTP, 31(356), 1087-1081.
2. S. Shima, M. Yang, 1992, Development of intelligent bending system, Form Tech Rev., 2(1), 63-72
3. Y. Saotome et al., 1992, Flexible V-bending system with simulation data-base and adaptive control, Proc. 43rd Japanese Joint Conf. for the Tech. of Plasticity, 409-412.
4. M. Yang et al., 1995, Development of intelligent V-bending system using data base (ibid.), Proc. 1995 Japanese Spring Conf. for Technology of Plasticity, 45-46.
5. T. Anzai, J. Endo, T. Mizuno and H. Yamada, 1997, Intelligent bending by controlling bending force, Proc. IPMM'97, (2), 1111-1117.
6. T. Anzai, J. Endou, et al., 1998, High accuracy V-bending by uniform force controlling (in Japanese), J. JSTP, 39(445), 70-74.
7. T. Otani, T. Oenoki, K. oda and M. Takada, 1996, Development of high accurate gauging in V-bending process with bent angle sensor, Proc. 47th, Japanese Joint Conf. for the Tech. Of Plasticity, 411-412.
8. K. Manabe, et al., 1988, Development of intelligent deep drawing system (ibid.), 29(330), 740-748.
9. K. Kawai, 1992, Development of spinning data base (ibid.), Form Tech Review, 2(1), 44-55.

Design of Enterprise Network Communication Subsystems

Adam Grzech

Institute of Control and Systems Engineering, Wroclaw University of Technology
Wybrzeze Wyspianskiego 27, 50-370 Wroclaw, Poland
Email: grzech@ists.pwr.wroc.pl

ABSTRACT

The aim of this paper is to present some basic rules to compare communication subsystems used in practice. The proposed procedures -- which can be treated as normalisation of practical approaches -- may be incorporated in design processes using expert systems to support the design of communication subsystems of distributed, corporate computer systems. Attention is given to formalise the simplest design task i.e., matching the corporate network communication subsystem to the requirements of each end-user system. The comparison of available and required subsystems is described using a predefined set of features at selected levels of granulation and detail.

INTRODUCTION

Theories of computer network design are rarely applicable in practice and so, development of non-standard configuration procedures have resulted. The field of artificial intelligence provides ways to design computer systems that work. Handling such problems relies on obtaining and formalising information, not only from theoretical knowledge, but also from design expertise. Such approaches have proven successful in creating practical rules used by experts for computer network design; in formalising the logic to solve computer network design problems; and in choosing the most suitable initial solutions for certain networking requirements. Requirement analysis is the first and most important step both in network design or tuning, but also in migration or integration tasks. No single aspect of the task is more important than a full understanding of the users' needs, since ultimately, they dictate the technologies, techniques and resources.

The proposed approach is a structured and adaptive technique that requires addressing the design problems of corporate network communication subsystems. It is based on the assumption that information processes and user requirements can be transformed into communication subsystem characteristics and vice versa. When known, this will allow us to compare categorised requirements (user's point of view) and available services (communication subsystem's viewpoint).

The problems of computer network classification, design, evaluation, planning, monitoring and management as well as, tuning, integration and migration tasks may be formulated in two different manners. The first is to find a communication system that is most suitable for a given data processing system while the second is to select a data processing system that exploits a known communication system. The first formulation is more common in practice and follows a top-down design methodology. Task formulation assumes that a set of requirement types, representative of a particular set of users, can be matched to the amount of offered services delivered by some providers and to select a provider offering an optimum level of services. This means that we must find a common way to express both user requirements and offered services.

ENTERPRISE NETWORKS

In a broad sense -- enough to encompass the different forms of past, present and future solutions -- the enterprise network must be defined as a corporate-wide network that ties together the processing, communications and storage resources of the corporation, thereby making such resources available to users distributed throughout the corporation. A more limited definition is focussed on the enterprise network in its

currently common form as a local area networks inter-network [4,6]. The generic architecture of enterprise network contains – at the most general level of generalisation - the following components:

- processors being able to run users programs and communicate using network protocols; processors are the sources and destinations of communications,
- distributed processing applications structured to allow different processes relationships,
- local communication between the processors on a peer-to-peer or client/server basis,
- premises-wide communication being an extension of various form of local communication,
- enterprise-wide communication being a collection of interconnected networks,
- interpremises network to assure off-net communication rather than within the bounds of a corporation,
- network management system for monitoring, billing, reconfiguration, security, etc. purposes.

Analysis, design and maintenance of enterprise networks do not focus on the transport of data. It is a broad field encompassing the creation, development and deployment of an infrastructure that comprehends user applications; network architectures, public standards, services, protocols, data transport techniques, etc. To meet requirements of contemporary enterprise networks it is necessary to obtain the optimal design. The latter covers various aspects of networking such as: network design, expansion, integration and migration planning, scheduling of implementation and resources allocation, prediction of services degradation and traffic overload, fault detection and isolation, network management, load balancing, etc. It is available only by introducing both theoretical knowledge and design expertise since no standard procedures for overall network design exists. Theoretically based design methods/algorithms are only minimally applicable in practise; they are isolated and devoted to well-defined and very limited in functional scope problems [1].

NETWORK DESIGN

Historically, network design has consisted of engineering primarily for the transport of voice or low-speed data traffic over dedicated private line facilities. The digital technologies have both increased network performance and enabled construction of huge networks capable of providing subscriber service for all forms of voice, data and image traffic and even, to assure specific quality-of-service parameters. The increased capabilities and performance of networks require more structured and adaptive techniques addressing the various aspects of network design problems. Rapid emergence of new technologies and the availability of a large variety of network equipment with different capabilities and different levels of technologies add a new range of variables to design consideration [6,7].

Network design can be viewed as a series of events from determining the basic requirements to managing the network after it is built. The network designer strives to achieve a balanced solution that comes as close as possible to meeting both the user and organisational demands. These needs are often in conflict: reaching the balance point between each end of the spectrum is the essence of network design influenced by business, technical, economical and political factors. The often-asked questions during the design process address various aspects of networking. The obtained solution may be selected to satisfy only present needs or may be designed to be ready to assure services for forecasted, future requirements. The network may be over-engineered or characterised by a large return of investment. The solution may be considered as application specific or as flexible and universal. The obtained solution may be within technical constraints or within corporate business requirements, etc.

The two major views of requirements are those of the user and those of the designer. The user looks at the network from the outside in, while the designer looks at it from the inside out, thus creating two myopic views which must merge to provide a comprehensive, complementary analysis beyond simple network ingress and egress design [4,6].

DESIGN TASK FORMULATION

Partitioning of a computer network is a simplification, but hierarchical decomposition is a major tool for complexity reduction. It is used precisely with the purpose to simplify the model of reality. As long as simplification is taken into account, hierarchical decomposition allows the focus to be on the system

components and their interrelations. The division of a network into distinctive and complementary subsystems may be considered based on hierarchical, multilevel nested computer and communication architectures. The subsystems co-operate and satisfy predefined requirements; the quality of delivered services depends on which levels are matched. At the design stage, the subsystems described by various sets of features must be transferred into descriptive forms that allow comparison. For the purpose of this paper, it is assumed that the design task may be formulated in two different manners. The first is to find a communication system, which is most suitable for a given data processing system while the second is to select a data processing system exploiting a known communication system. The first formulation is more common in practise and follows the top-down design methodology [2,6].

The quality of the design process depends on a set of factors taken into account as well as procedures that establish the parameters of the considered systems. The process defining the set of factors taken into account at the design stage is based on identification and recognition of issues providing strong characterisation of the considered data processing system and communication subsystem. The aims of the identification and recognition stages are to extract characteristics that well-define (for purpose of further design) as well as, properly recognise the considered processing and communication systems among other systems requiring and delivering a qualitatively and quantitatively comparable scope of services. The aim of procedures to compare the considered systems is to adopt some assumed hierarchy of facts that reflect collected knowledge and experience which influence the quality of the design process [2,3,5].

Different design task formulations assume it is possible to match requirements (representing particular data processing subsystems) to the amount and quality of services offered. Moreover, it is assumed that some performance measures when optimised, allow us to select the best interrelations. This means that a common way to express both required and delivered services is necessary. Assumed description of the data processing subsystem should lead to a set of observable and measurable parameters translated onto values of parameters that define the communication subsystems [3].

DATA PROCESSING SUBSYSTEM DESCRIPTION

Let us assume that the architecture of any data processing subsystem is considered, for the purpose of design, as a hierarchy of L distinguished, layered and nested components (mechanisms, procedures, algorithms, etc. distinguished at the assumed level of granulation). The set of all components is denoted by Ω . Each l -th component ($l=1,2,\dots,L$) represents some unique functionality within the general data processing architecture. Localisation of the l -th component in the overall hierarchy is determined by its functionality. This is equivalent to scoping the services offered to other components located at the $(l+1)$ -th and above, as well as, scoping services required from components located at the $(l-1)$ -th layer and below.

The functionality of the considered l -th component is understood as assuring some data processing subsystem feature. The set of all values of the l -th feature is denoted by $\Omega_l = \{\alpha_{li}\}$ where α_{li} is an i -th value of the l -th feature, where $i=1,2,\dots,I_l$ and I_l is the number of all distinguished values of the l -th feature. For purpose of simplicity, it is assumed that l -th component functionality is equivalent to a unique feature. It is worth noting that the feature values may be numbers (quantitative measures) or descriptions (qualitative measures). Knowledge of Ω is equivalent to possessing subjective, general-purpose language allowing common descriptions of a variety of data processing systems at the assumed level of granulation. The applicability of the set Ω may be reduced to some classes of data processing subsystems achieved by increasing the granulation of description may encompass the reduction.

The Ω set may be obtained during process observation when real design tasks are performed. Let us assume that any real data processing subsystem is the result of co-operation among many different, existing components forming a set denoted by Ω_R . For design purposes, some components belonging to set Ω_R are neglected and discarded. The set of neglected components may include -- among others -- components having no impact on the considered data processing environment design or components characterised by unknown, non-observable or non-measurable influences on a designed network performance. Reduction of the set Ω_R leads to a set denoted by Ω_P ($\Omega_P \subset \Omega_R$). It contains all components to be treated as

objectively important for design process purposes. Further analysis (based on theoretical and/or practical knowledge) of components belonging to the set Ω_P leads to its reduction to set Ω_D ($\Omega_D \subset \Omega_P$) containing a set of recognised features -- in contrast to Ω_P -- subjectively important for the design process at the assumed level of granulation. The relation $\Omega \subset \Omega_D$ reflects the fact that the same set of components within various hierarchies may produce different design results; Ω refers to a set of components as well as a particular hierarchy of the selected components. Changes in sets of components and hierarchy of the components lead to another description language of data processing systems.

DATA PROCESSING SUBSYSTEM CLASSIFICATION

By treating the Ω set as a common language, we can distinguish classes of data processing subsystems; the set of all classes is denoted as $D(\Omega)$ where $D(\Omega) = \{D_1(\Omega), D_2(\Omega), \dots, D_K(\Omega)\}$. This notation reflects the fact that different sets of components within different hierarchies lead -- in general -- to different sets of classes and to different descriptions of these classes. The set $D_k(\Omega)$ ($D_k(\Omega) \subset \Omega$ for $k = 1, 2, \dots, K$) is composed at most of L components ($D_{kl}(\Omega)$, $l = 1, 2, \dots, L$) equal to the number of components of the set Ω . The set $D_{kl}(\Omega)$ ($D_{kl}(\Omega) \supset d_{klj}$ where $j = 1, 2, \dots, J_l$) if it exists contains values of the l -th feature from the set $\Omega_l = \{\alpha_{li}\}$, i.e., $d_{klj} \in \{\alpha_{li}\}$, $J_l \leq I_l$. The k -th class contains solutions denoted by $D_k^{(u_k)}(\Omega)$ where $u_k = 1, 2, \dots, U_k$ and U_k is the number of all the different subsystems within the $D_k(\Omega)$ class. In general, the following conditions are satisfied:

$$\begin{aligned} \forall_{k \in \{1, 2, \dots, K\}} \exists_{l \in \{1, 2, \dots, L\}} D_{kl}(\Omega) &\neq \emptyset, \\ \forall_{k \in \{1, 2, \dots, K\}} D_{kl}(\Omega) &= \Omega_l \text{ or } D_{kl}(\Omega) \subset \Omega_l \text{ or } D_{kl}(\Omega) = \emptyset, \\ \forall_{k, j \in \{1, 2, \dots, K\}, k \neq j} D_k(\Omega) &\neq D_j(\Omega) \text{ iff } \exists_{l \in \{1, 2, \dots, L\}} D_{kl}(\Omega) \cap D_{jl}(\Omega) \neq \emptyset. \end{aligned}$$

This expression means that at least one feature is necessary to define the class of a particular data processing subsystem and that any two classes differ at least by one value of any one feature.

COMMUNICATION SUBSYSTEM DESCRIPTION

Let us assume -- similarly as in the previous section -- that architecture of any communication subsystem is given as a hierarchy of N distinguished, layered and nested components (mechanisms, procedures, algorithms, transmission media, etc. distinguished at an assumed level of granulation). The set of all components is denoted by Σ . Each n -th component ($n = 1, 2, \dots, N$) represents some unique functionality within general communication subsystem architecture providing some communication subsystem feature.

The role of the n -th component in the assumed hierarchy is determined by the scope of functionality and co-operating components from the higher and lower layers. The set of all values of the n -th feature is denoted by $\Sigma_n = \{\beta_{nm}\}$ where β_{nm} is an m -th value of the n -th feature, $m = 1, 2, \dots, M_n$ and M_n is the number of all distinguished values of the n -th feature (the assumption that the n -th component functionality is equivalent to a unique feature may be easily generalised). The Σ set contains a set of features recognised (from designer theoretical and/or practical knowledge) as subjectively important for the design process purposes at the assumed level of granulation. Selection of the set Σ means selection of some set of components, as well as, a particular hierarchy of the selected components uniquely determining communication subsystem services and service features.

COMMUNICATION SUBSYSTEM CLASSIFICATION

The Σ set allows us to distinguish classes of communication subsystems; the set of all classes is denoted as $C(\Sigma)$ where $C(\Sigma) = \{C_1(\Sigma), C_2(\Sigma), \dots, C_S(\Sigma)\}$. In general, different sets of components and different hierarchies lead to different sets of classes and to different descriptions of each class. The set $C_s(\Sigma)$ ($C_s(\Sigma) \subset \Sigma$ for $s = 1, 2, \dots, S$) is composed at most of N components ($C_{sn}(\Sigma)$, $n = 1, 2, \dots, N$). The set

$C_{sn}(\Sigma) (C_{sn}(\Sigma) \supset c_{snt}$ where $t = 1, 2, \dots, T_n$) may contain none, some or all of the values of the n -th feature from the set $\Sigma_n = \{\beta_{nm}\}$, i.e., $c_{snt} \in \{\beta_{nm}\}$, $T_n \leq M_n$. The set $C_n(\Sigma)$ contains solutions denoted by $C_s^{(w_s)}(\Sigma)$ where $w_n = 1, 2, \dots, W_n$ and W_n is the number of all different subsystems within the n -th class. In general, the following conditions are satisfied:

$$\begin{aligned} \forall_{s \in \{1, 2, \dots, S\}} \exists_{n \in \{1, 2, \dots, N\}} C_{sn}(\Sigma) &\neq \emptyset, \\ \forall_{s \in \{1, 2, \dots, S\}} C_{sn}(\Sigma) = \Sigma_n &\text{ or } C_{sn}(\Sigma) \subset \Sigma_n \text{ or } C_{sn}(\Sigma) = \emptyset, \\ \forall_{s, t \in \{1, 2, \dots, S\}, s \neq t} C_s(\Sigma) &\neq C_t(\Sigma) \text{ iff } \exists_{n \in \{1, 2, \dots, N\}} C_{sn}(\Sigma) \cap C_{tn}(\Sigma) \neq \emptyset. \end{aligned}$$

This expression means that at least one feature should be available to define the classes of communication subsystems in which any two classes differ by at least one value of any one feature.

COMMUNICATION SUBSYSTEM DESIGN

The formulation and solution of various design tasks in the area of computer-based systems using traditional analytical and simulation tools are only minimally effective. Applicability of such methods is limited to the consideration of the need to set up design process descriptions with granulations different from those which are used or available in computer-based and telecommunication system theories. There is a lack of standard computer system configurations, so that -- in most cases -- the computer-based system designer is the ultimate authority. This means that practical design of such systems is based in great measure on experience obtained in areas of data collecting, and in theoretical and practical knowledge formalisation. Direct comparison of features and feature-values describing data processing and communication subsystems is practically impossible. It is mainly because these two system classes are completely different in nature and are characterised by applying different qualitative and quantitative variables and parameters. The difficulty in direct comparison of these subsystems gives the reason that we must apply various procedures leading to the design solution. Comparison can be done by applying various individual and more-general procedures.

To solve the design task, i.e., to select the proper subsystem based on comparison of the required/offered (data processing) and offered/required (communication) subsystem services, descriptions of both subsystems, as well as, procedures to translate and standardise the descriptions in both directions are necessary. For the purposes of this paper it is assumed that the design process is devoted to selecting a communication subsystem from known data processing subsystems is more-natural and more-frequently performed in practice. Two different design strategies may be distinguished: direct and indirect.

Qualitative and quantitative features applied to describe data processing and communication subsystems, in general, have different physical natures, dimensions and priorities, as well as, being characterized by different measures of quantities and different languages to provide descriptive qualities. To solve the design task, i.e., to perform classification, comparison and optimisation of subsystems belonging to two different spaces must be standardised and transformed; especially when numerical values are required.

Standardisation and transformation may be understood as processes that lead to numerical representations of various subsystems. Let us assume that numerical representations of data processing $D_k^{(u_k)}(\Omega)$ and communication $C_s^{(w_s)}(\Sigma)$ subsystems are denoted by $x_k^{(u_k)}$ and $y_s^{(w_s)}$ where $x_k^{(u_k)} = TR(D_k^{(u_k)}(\Omega))$ and $y_s^{(w_s)} = TR(C_s^{(w_s)}(\Sigma))$. The elements of vectors $x_k^{(u_k)}$ ($x_{kz}^{(u_k)}$, $z = 1, 2, \dots, Z$) and $y_s^{(w_s)}$ ($y_{sy}^{(w_s)}$, $y = 1, 2, \dots, Y$) are computed from a set of rules denoted as $\omega_{kz}(\dots, d_{k1j}^{(u_k)}, \dots, d_{kLj}^{(u_k)}, \dots, d_{kLj}^{(u_k)}, \dots)$ and $\omega_{sy}(\dots, c_{s1t}^{(w_s)}, \dots, c_{snt}^{(w_s)}, \dots, c_{sNt}^{(w_s)}, \dots)$ respectively for data processing and communication subsystems. These rules express the influence of various features values on the value of particular elements of the above vectors. The rules are based on assumed, subjectively-defined relations and reflect the importance of features and feature-values defined by the designer at the formulation design stage. In the above expressions, the numbers Z and Y denote dimensions of vectors x and y . For the purpose of subsystem

comparison, the dimensions of the vectors attached to all solution should be the same within the various classes, i.e., $D(\Omega)$ and $C(\Sigma)$ are the same.

The first, direct, design strategy is a one-step optimisation procedure. The aim of this optimisation is to find for a given data processing subsystem D , a communication subsystem $C_I(D)$ belonging to one and only one class $C_s(\Sigma)$ ($s = 1, 2, \dots, S$) from the $C(\Sigma)$ set. This is equivalent to finding s° and w_s° such that:

$$C_{s^\circ}^{(w_s^\circ)}(\Sigma) = \min_{s \in \{1, 2, \dots, S\}, w_s \in \{1, 2, \dots, W_s\}} f(TR(D), TR(C_s^{(w_s)})),$$

where $f(.,.)$ measures the "distance" between transformed and standardised descriptions $TR(D)$ and $TR(C_s^{w_s})$ of a data processing subsystem D and various $C_s^{w_s}$ communication subsystems, respectively.

The second, indirect, design strategy is based on a two-step procedure. At the first step (classification step), the particular data processing subsystem D is classified as belong to one class $D_k(\Omega)$ ($k = 1, 2, \dots, K$) from previously established sets of data processing subsystem classes. The aim of the second step (optimisation step) is to determine the solution $C_{II}(D)$ from one and only one class $C_s(\Sigma)$ ($s = 1, 2, \dots, S$). The task is equivalent to finding two pairs: (k^*, u_k^*) and (s^*, w_s^*) such that:

$$D_{k^*}^{(u_k^*)}(\Omega) = \min_{k \in \{1, 2, \dots, K\}, u_k \in \{1, 2, \dots, U_k\}} g(TR(D), TR(D_k^{(u_k)})) \text{ and}$$

$$C_{s^*}^{(w_s^*)}(\Sigma) = \min_{s \in \{1, 2, \dots, S\}, w_s \in \{1, 2, \dots, W_s\}} h(TR(D_{k^*}^{(u_k^*)}), TR(C_s^{(w_s)}))$$

where $TR(D)$, $TR(D_k^{u_k})$, $TR(D_{k^*}^{u_k^*})$ and $TR(C_s^{w_s})$ are transformed and standardised descriptions of data processing and communication subsystems respectively while $g(.,.)$ and $h(.,.)$ are performance measures applied to solve the classification and optimisation tasks, respectively. Solutions obtained from the two strategies and given by two pairs (s°, w_s°) and (s^*, w_s^*) are generally different.

CONCLUSION

One of the most important research motivations behind implementing the proposed procedures is their ability to behave exactly like an individual designer. This feature is obtained by providing the possibility to select the set of factors and to propose their hierarchy as being important to design process results. The method can give a solution, answer questions that fall within the domain of the considered subsystem and reason solution results by employing various methodologies, knowledge and experience representations as well as different predefined or ad hoc applied tools in an individually pre-established order.

REFERENCES

1. H.I. Fahmy, Ch. Douligeris, 1995, END: An expert network designer, IEEE Network, 12, 18–27.
2. A. Grzech, 1997, Design of communication subsystems of computer networks, Proceedings of the 3rd Conference on Knowledge Engineering and Expert Systems, 2, 69–77.
3. A. Grzech, 1998, Qualitative and quantitative comparison of communication subsystems of computer networks, Proceedings of the International Symposium on Systems Modelling Control, 346–353.
4. R.A. Mercer, 1996, Overview of enterprise network developments, IEEE Communications Magazine, 34(1), 30–37.
5. K. Nowak, 1998, Design the distributed computer systems supporting co-operative working, Proceedings of the 13th International Conference on Systems Science, 3, 139–150.
6. D.L. Spohn, 1997, Data Network Design, McGraw-Hill.
7. R.J. Wieringa, 1996, Requirements engineering, John Wiley & Sons.

The Industrial Desktop – Real Time Business and Process Analysis to Increase Productivity in Industrial Plants

Osvaldo A. Bascur

OSI Software, Inc.
Houston, Texas, USA

ABSTRACT

Intelligent systems (IS) technologies have received much attention in a wide range of process engineering applications including process operations. Rapid change in applying the latest technologies has become a serious challenge to both management and technical teams. Objects and components are changing the way we relate to our computer and networks and most feel the rate of change will continue to increase.

All information technology systems have data and communications tools for personnel. The industrial desktop can be adapted to automate decisions, to intelligently analyze large amounts of data, and to learn from past experiences whether from operators, engineers or managers. They can adapt their desktop according to their domain knowledge, roles, skills and responsibilities. Collaboration between functions (Operations, Management, Maintenance, and Engineering) is enhanced. Access to data and analysis tools enables plant personnel to try new ideas, determine and track the right targets, determine and track the best patterns, and transform and store data and knowledge.

This paper presents a description of the data hierarchy and analysis means needed to improve process operations. Continuous improvement with an innovation loop fueled by data collection and analysis methods emerges as the best method for active decision-making and collaboration. A sampling of results include extending sub-critical equipment availability, increased production by faster detection of process bottlenecks and operating cost reduction. Descriptions of these applications are presented.

CHALLENGES IN INSTRUMENTATION, CONTROLS AND MANAGEMENT

Process control and software developments in the process industries over the past decade have been greatly influenced by new instrumentation advances, size/speed of microprocessors and software developments. These developments provide opportunities for the plant personnel to effect strategies that permit them to operate the process units more profitably. In spite of these advances, a study in the pulp and paper sector by Entech Control reports that 80% of distributed control system loops actually increase process variability as compared with manual control systems [1]. The main reasons for variability were:

- 20% due to design causes
- 30% due to control tuning
- 30% related to equipment performance.

In order to maintain greater up-time for any strategy, it is necessary to measure performance and continuously improve behavior [2]. The most effective tool is a continuous, real time, on-line audit called a control monitor. Every day, it will rank the loops and verify that the business needs are still valid and are being met. This list is used to adjust instrument maintenance, training of the operators and a program for management of change for the controllers. At one paper mill, they reported that monitoring plus routine action were helped them increase up time of critical controls from less than 50% to over 80%.

An emerging industrial plant index called the overall process and equipment effectiveness is given by the plant availability (A) times the performance efficiency (PE) and times the rate of quality (RoQ); i.e.,

$$OPE = A \times PE \times RoQ \quad \text{where,}$$

A = Availability or unscheduled downtime plus scheduled downtime

PE = Performance Efficiency (PE) or idling and minor stoppages, reduced speed of equipment

RoQ = Rate of Quality (RoQ) or rework, yield or recovery loss.

This expression gives us the net operating time, limited by minor stoppage and speed losses. It also considers quality defect and yield losses. The variable operating time equals the operating time minus downtime losses, speed losses and quality losses. This same approach works with all types of equipment, not just controls. Equipment problems need to be identified and analyzed quickly; only then can they be solved. Most companies already have excessive data about their equipment stored in various locations:

- Maintenance Management Systems
- Cost Systems
- Predictive Systems
- Production Systems
- Manufacturer specifications and reliability data

What is needed is an environment that simplifies integration of the data with the analysis tools. We must also prepare for the fact that with the new communication media, the Internet, field buses and intelligent assets (e.g. motors), there will be orders of magnitude per year increases in the amount of data available. People in the plants need tools to simplify decision-making, intelligently analyze large amount of data, and learn from their mistakes [3].

The data hierarchy framework and communication infrastructure must be restructured to improve analytical capabilities by utilizing the software available in plant PC desktops. Most measurement problems arise from lack of data for analysis or adequate analysis tools. The epoch development tool is the graphical user interface with the ability to embed for analyzing, accessing, discovering and taking actions from the plethora of data from a plant. With components, this is all accomplished with reusable code that promotes continuous improvement. Proper design of this interface makes collaboration between operations, engineering and management possible and necessary [4,5,6,7,8].

The Internet adds a new requirement. Many users who share information and analysis methods with plant personnel are located on the company Intranet and use a browser to view plant operation and make their decisions. This means these, or similar, components must be deployable by an Internet Server into a browser - at least in a limited mode. Microsoft has introduced their DNA (Digital Internet Architecture) for Manufacturing based on the idea that the Internet is the ultimate process re-engineering tool.

THE INDUSTRIAL COMPONENT DESKTOP

All business processes and the software that supports them must incorporate reusable components and promote continuous improvement. This is a key improvement for traditional plant Process Information Systems (PIMS). A traditional PIMS is open-loop. Information is gathered into databases and disseminated as reports and on-line inquiries to all requesters; here the system's responsibilities end. This system does not promote sharing information or transforming it to knowledge **for action**. It must have a development framework to support the innovations needed for a more aggressive approach.

The new Microsoft Windows 2000 and Office 2000 software that runs on a plant PC, called the industrial desktop, is central to the development environment needed to support innovation. It allows the user to increase his or her effectiveness by adding components to their familiar environment including video views, telephones, search algorithms, or even, custom components created with user level tools. This changes the individual from a fire fighter to a proactive worker who can analyze, make discoveries about the plant and business processes and, most importantly, implement his/her findings. The user tailors his/her environment according to his/her role, skills, responsibilities, and accountability in the plant system.

Why should they do it themselves? The main reason is that all of these changes and ideas are also very difficult and expensive to transfer from one person to another and, when others do to the work, much of the gratification is lost and a training issue arises. In the past, these little mini-projects seldom get done. The rewards of implementing these small ideas are large and as an aggregate they become a large, high-benefit change [9].

INTEGRATED DATA, ACTION, and HUMAN INTERFACE SYSTEM STRUCTURE

A Generic Data Structure for Process Control (GDSPC) was constructed [10] to review the different tools available for process management and control in the process industries. A schematic representation of the GDSPC is shown in Figure 1. This framework describes how data are transformed into information and how the different levels of information for the data, action and user interface are related. The information is processed by modules at each level in a timely manner. All information is saved for historical analysis.

The data processing hierarchy is composed of:

- The data acquisition subsystem,
- The data validation subsystem,
- The data classification and analysis subsystem,
- The data verification and estimation subsystem,
- The plant data management subsystem and,
- The plant data optimization subsystem.

The first four levels of the data processing are real time tasks. They need to be processed on a time class demand. Plant data management and plant data optimization processing are performed depending upon operation planning and economic evaluation time frames.

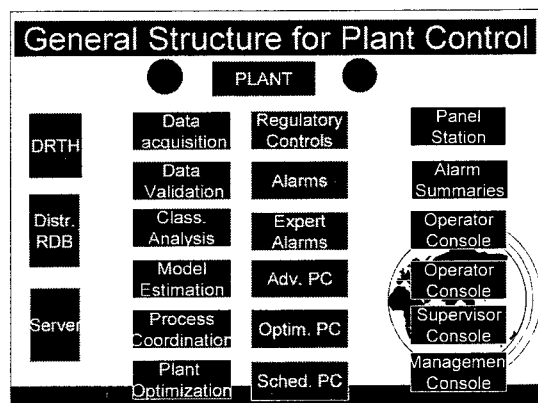


Fig. 1. A modular approach to manage data, actions and analysis.

The Data Acquisition Subsystem

The first layer is the acquisition of measurable process variables from transducers and output of variables to the actuators, and the digital input and output. At this level, instrument diagnosis can be implemented. Even with the use of smart transmitter, these cannot detect, not correct for, a drift or span error until the transmitter output has been compared with known signals or values at operating conditions [11].

The Data Validation Subsystem

The data validation step can be implemented by several methods. The traditional method is to check for validity by using a known relation that must hold true. Each process measurement and resultant calculation is compared against high and low limits as well as a maximum rate of change. Depending on the type of process, a voting technique can be used when the measurement is so critical that two sensors can be economically justified. The data verification results are posted to an alarm summary and/or a data logger.

The Data Classification and Analysis Subsystem

The classification and analysis of validated process data divides the process into several modes or operational status. This classification is facilitated by incorporation of common-sense rules dictated by expert operators and process engineers. These triggers are calculated indexes, which define a state of the process unit or equipment to define a production, quality, cost, environmental or equipment status.

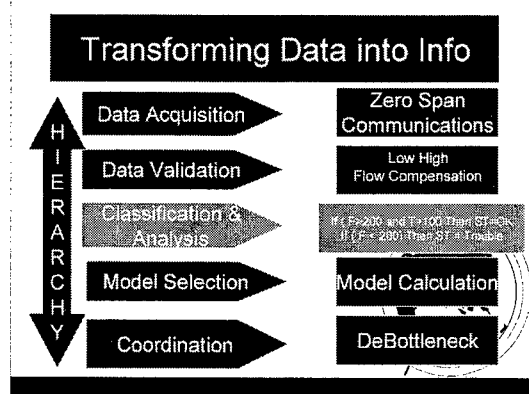


Fig. 2. Plant data hierarchy example.

The trigger indicator is used in real-time applications or to enhance historical data extraction for calculation of mass balances, total operating hours under certain conditions, or to select the type of model to use in inferential or economic functions.

The utilization of several timeframes of a data point such as data average, moving-average, minimum, maximum, standard deviation and rate of change can be utilized to classify the faults or process status. This history of process values, events and alarms is usually available but rarely automated. These time-series of variables are available to extract pattern features to check for deviations from normality. A node-type system is used to minimize the computer processing time. In this generic-data-action hierarchy, a human-equivalent metaphor is used. For example, if the communications is broken, or if the sensor is in trouble, or if the process is in a certain status, special detailed detection is set for further data processing.

Communication status node - the first item to check is the communication link between control subsystems. Any missing process information required by the unit advisor must be detected and acknowledged as soon as possible. Detection of on-line analyzer malfunctions must be included in this data analysis subsystem.

Equipment monitoring and diagnosis (Trouble node) - the data subsystem should recognize the unit equipment sensors and actuator limitations and the safety restrictions. The data trouble subsystem should act as an advisor to ensure that process-states do not violate such limits. The history generated by these types of points is used for predictive maintenance of instruments and sub-critical equipment. Technical methods to detect and respond to failures are documented in the literature [12,13], Diagnostics and Reliability-Based Maintenance, <http://www.rosemount.com/products/ams/rb-maint.htm> (Internet), May 8, 1998. In addition, simple checks can be implemented such as power-draught and oil pressure in rotating machines, vibration and amperage in fast rotating equipment, pump efficiency estimation, rate of change to detect leaks in fluid lines, etc

Process operational status node - this data subsystem helps an operator by advising him of tasks requiring manual intervention. The operator will only take care of process units that require special attention and will try to solve plant bottlenecks. Data classification and analysis should prevent over-controlling process units. In this manner, the operator can have a global vision of the situation and be prepared to solve real trouble. The availability of correct action response can have a large effect on plant performance.

Control advisor node - depending on the type of economic objective there are several forms of control available for a process unit. For example, when a preprocessing part of the plant is down, the next subprocess may slowdown and change from a maximizing throughput objective to an improved product quality objective. This control strategy will be set in time following a guided procedure with proper interlocks on the control stations. For example: a slowly varying system requires analysis of historical data acquired over several hours.

The Data Verification and Estimation Subsystem

After data has been classified and process operational status identified, further process inference can be obtained by using estimates of unmeasured variables for implementation of advanced control strategies. If the process is considered to be at steady state, a form of data reconciliation can be used. A mass-balance closure adjustment can be used. The adjustments are based on a weighted-least-squares fit subjected to the mass flow rate and composition mass balance techniques.

The simpler techniques use multiple linear regression techniques to infer the property from other variables. Sometimes stochastic features are added to include plant noise. In other instances thermodynamic models and process measurements are used to infer properties. (End-point, freeze point, flash point, cloud point, calculations in fractionation systems, etc.).

A state variable approach can also be used to provide optimal estimates using Kalman-filtering techniques [14]. In this method, the estimates are calculated by proper weighting of information provided by a process model and the measured variables. The weighting is based on statistical properties of the process and the available measurements. These techniques add a great value to current instrumentation with a minor investment in proven software for an existing distributed control system. Estimation techniques can be used very effectively when an on-line analyzer is not justifiable or reliable. The process visibility obtained by these methods allows better design of the controllers and plant accounting and planning.

The Plant Data Management Subsystem

Once the missing data for the process units is acquired the coordination of these units is possible. At this point each process unit is controlled in a sub-optimal fashion if other unit interactions are not taken into consideration. To provide for a sound operation the coordination of the process units needs to be incorporated. Items such as production rates, maintenance schedules, energy considerations, capacity limitations, demand and process constraints are set according to plant economic objectives.

The presentation of the plant data in adequate reports improves for further decision making. Plants, sections or units material and energy balance reports, calculation of global performance indexes, inventory of raw materials and products, etc. can be facilitated by a reliable distributed data base system. The availability of special functions (totals, averages, searches, statistical analysis, time calculations, etc.) increases the efficiency of data consolidation and presentation.

Traditionally, supervisors set the operating conditions and communicate to their operators. The industrial desktop enables closing the loop between operations and management. The data coordination level is also valid by defining the level of collaboration between operation and maintenance functions. Implementation of these loops avoids chasing of conditions, which sometimes might result in instability.

The Plant Data Optimization Subsystem

At the top of the hierarchy the data obtained is processed to evaluate the global profitability of the operation. The validated data and calculated technical indexes should be combined with operational costs and production indexes. The data forecast resulting from consideration of the current market conditions and plant availability are fed to the planning functions to set the maximum profitability. At this level, the development at the industrial desktop towards planning, act and measure are the key for achieving high results. The integration of process efficiency and equipment operability functions with unit availability status and anticipated scheduling commitments are defined. Data validation, classification and verification are vital since the optimum will depend on having accurate data and a validated process model.

The data is used to develop a recommended action plan for plant operations to balance limited resources, operating cost variances, estimated repair costs, utility costs, and projected revenue opportunities.

PROCESS ANALYSIS AND DEVELOPMENT

Emerging new data analysis technologies are becoming available to expand the detection of cause and effects of multivariate systems. Additional tools are emerging to facilitate analysis of multivariate systems such principal component analysis PCA.

When processes are subject to large unknown disturbances, the alternative of applying operator mimic as been successfully used [2,15]. A degree of imprecision can be assigned to the generated procedures, which can be exercised under a real time environment. The use of time derived variables are highlighted in the development of a Process Control Matrix and the Construction of Historical based

Process Knowledge Table				
Manipulated Variables	Controlled Variables			
	T Chamber	Chiller	PLL	Press
Gas #1	Fast Down	Slow Up	Slow Up	
Gas #2	Fast Down	Slow UP	Slow Up	
Gas #3				
Power	Fast Down			
Fpower				
Rpower				
Impedance	Up Down	Fast Down	Fast Up	

Fig. 3. Process knowledge table.

The concept of approximate reasoning can be used to model the process in question [16,17,18]. The representation of cause and effect of a process model is captured using a linguistic approach. Tables are generated using empirical or deterministic models. The key is that the model has a degree of fuzziness. The process knowledge table describes the effect of control variables by manipulated variables subjected to the observed variables (Figure 3). It uses a fuzzy description of the response or historical pattern. The observed variables are usually slow moving variables which define the operating conditions such as the age of a well, age of the catalyst, pump characteristics, equipment wear, or the cleanliness of the reactor.

The decision knowledge table also uses linguistic approximations and developed using time-derived variables (Figure 4). These auxiliary variables are used to develop the trigger points necessary to develop the operating condition status indices. The use of the historical data is stressed to define at all times the process pattern. The subset of data used to develop the soft limits can be obtained by classification of the variables for each of the triggered indices.

Decision Table				
Measured Variable	Snap Shot	Mov-Avg (time)	Rate (Time) (Time)Avg Std	(e)2 (time)
Temp1	high	high	hrate med mx	min
Flow1	H/L	H/L	low	
Level		low	hrate	min
Flow2		mhigh		
Pres2		mhigh		
Additional Step, Phase	xyz	SLOW MOVING FUZZY VARIABLES		
Age		low	old/mold/new/unknown	
Cleaning		low	recent/medium/close	
# comp.		high	few/many,very many	

Fig. 4. Process operating condition decision table

Figures 3 and 4 show an example of a process knowledge table generated using the presented approach. The trigger point calculation are implemented either for real-time execution at the server level for on-line diagnosis or at the client for process data analysis and development of the trigger actions.

PROCESS TROUBLESHOOTING EXAMPLES

A good example is the Microsoft interactive software agent (see Figure 5.). Basically, this agent provides a conversational user interface, which is employed to enhance, rather than replace, the usual Windows graphical interface. Activation of this agent can be set by a real-time event called a change event. As such, the Agent can explain why key process indices have gone off-spec. Activation of an agent can use a described method to define a trigger based on sharp or fuzzy rules. The Agents (Genie, Merlin and Robbie) have a palette of animations, a synthesized voice, the ability to synchronize their movements with recorded voices, and are able to respond to mouse clicks, internal events and voice commands. You can implement agents through a COM interface or ActiveX control. You can tailor an example to your own needs [19]. In Figure 5, the client software container has a VBA which enables automation of the genie based on an event generated on the desktop.

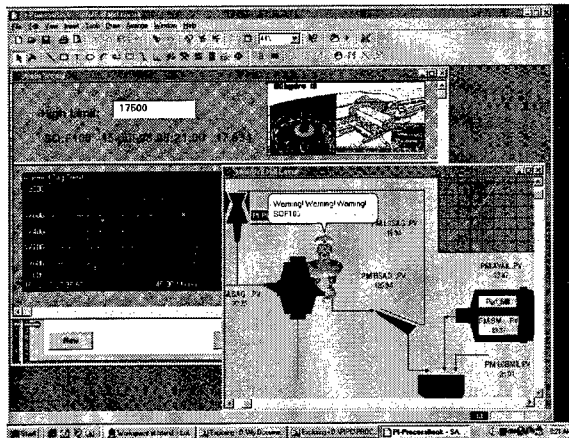


Fig. 5. Microsoft Agent triggered by a change event on the Industrial Desktop.

Spreadsheet Real Time Filtered Data Extraction and Exception Reporting

Instead of searching, and then, retyping information into the spreadsheet, Figure 6 shows a spreadsheet is set to run every morning a 8 AM and get the data and exceptions for a metallurgical temperature of a reactor. Every morning, filtering is executed prior to arriving at the office. (Note: This exercise can be scheduled at any desired time - minute, hour, shift, etc.). At the same time, correlation between the dependent and independent variables is performed to check relationships between the variables. If these correlation coefficients change, the relationship is no longer valid and adaptation of the model is required.

The extraction data requires special filtering to obtain a data subset. These filters act as special methods to reconstruct the most appropriate time-series. Usually, quality, process equipment, and environmental data are collected at different scan times than process data. Special methods are available to reconstruct the data set for scan times available for key indicators. It is important to note that once a spreadsheet template is built, it can be re-used for other time intervals. As such, it can be used for fault diagnosis at any time.

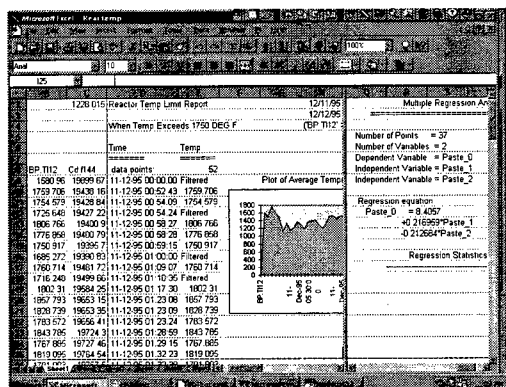


Fig. 6. Automated exception report and process data correlation.

Multivariate statistical methods are very effective at extracting hidden information in problems with multiple correlated variables. Several new data analysis methods have been described [20, 21, 22, 23, 24].

In addition, having simple methods to manage the process data enables to develop neurofuzzy applications [25]. Fuzzy tools can be added-in to Excel to generate rules to develop models for analysis of data. A data extractor such as PI-Datalink can be integrated by a fuzzy rule generator.

The workstation real-time environment described above can be enhanced by using mathematical objects available in the market. These software packages are usually data-intensive and they need data classification methods as discussed earlier.

Although the function here is quite simple, the underlying infrastructure necessary to make this happen is not simple. Effectively, the user can customize or manipulate PI objects in any manner desired. For example, if a user has a database of start-ups and tests, he/she could develop an application that automatically sets the time on the displays back to that time frame. This database could reside in the Batch tracking module or some external database.

Exemption Reporting in Batch Processes

Figure 7 shows a data analysis tool to detect if a batch in progress or done meet the specifications of the golden batch. Such technique permits to analyze batches without the need to wait for further processing or quality control. A subset of good batches is used to obtain a statistical representative high and low pattern limits for each of the phases in the batch process. Once, it is shown that if the key variable has violated more than a certain amount of times the acceptable envelope, the product can be discarded or recycled.

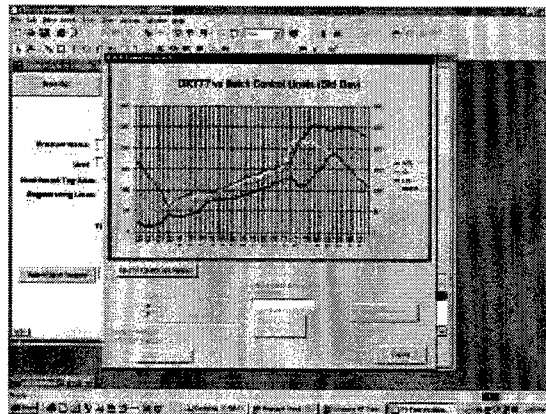


Fig. 7. A graphical representation of batch exemptions compared to the golden batch pattern.

CONCLUSIONS

A robust environment on the industrial desktop provides real time, historical process/equipment information and business information for all functions in the process industries. This environment enable users to build, construct and maintain their views of the plant for simplified performance monitoring, process and equipment troubleshooting, continuous improvement and innovation.

In general the results can be summarized as follows:

- Improved quality and speed of plant monitoring and diagnosis
- Optimizing plant efficiency performance (by continuous improvement and innovations)
- Condition based maintenance
- Improved production and regulatory information
- Consistent and comparable information across the enterprise (benchmarking).

These technologies are available today to rethink plant operations and to increase the performance of existent production systems. The key to re-engineering is linking people, business processes, strategies and the best enabling technologies. Many benefits are available that do not require disruptive re-engineering.

REFERENCES

1. Bialkowsky, 1996, Auditing and Reducing Process Variability in Your Plant, Presentation at Fisher Rosemount System Users Group, November, Houston.
2. Bascur, O.A., 1991. Human Factors and Aspects in Process Control Use, Plant Operators Symposium, SME, Colorado, 85-94.
3. Zuboff, Shoshana, 1991. Human Potential in the Age of the Smart Machine. The Consultant Forum, 6(1), 2-6.
4. Bascur, O.A., 1993, "Bridging the Gap Between Plant Management and Process Control," Emerging Computer Techniques for the Minerals Industry, B.J. Scheiner et.al, Eds. SME, Littleton, CO, 73-81.
5. Bascur, O.A., Vogus, C.B. and Bosler, W.H., 1992 "Long Term Knowledge Integration With OSHA PSM," National Petroleum Refinery Association, Computer Conf., Wash. D.C., November 16-18.
6. Kennedy, J.P., 1996. Building an Industrial Desktop, Chemical Engineering, January.
7. Bascur, O.A., Kennedy, J.P., 1995. Measuring, Maximizing and Managing Performance in Industrial Complexes, IMPC Proceeding, SME.
8. Bascur, O.A., Kennedy, J.P., 1996. Industrial Desktop -Information Technologies in Metallurgical Plants, Mining Engineering, September.
9. Bascur, O.A., Kennedy, J.P., 1999. Real Time Data Management to Increase Profitability of Mining and Metallurgical Operations, Conference of Metallurgists, Quebec City, August 21-25.
10. Bascur, O.A., 1988. A Control Data Framework with Distributed Intelligence, ISA Proc., Pap.88-1556.
11. Berge, J., 1998. Fieldbus Advances Diagnostics, In Tech, April, pp. 52-56.
12. Himmelblau, 1978. Fault Detection and Diagnosis in Chemical and Petroleum Processes, Elsevier Scientific Publishing Company.
13. Pau, L.F., 1981, Failure Diagnosis and Performance Monitoring, Dekker, New York.
14. Gelb, A. (Ed.)1974. Applied Optimal Estimation, The MIT Press, Cambridge, MA.
15. Bascur, O.A., 1990. Expert Process Advisor, Control '90, R.K. Rajamani, J. Herbst Eds., SME, 67-76.
16. Zadeh, L.A., 1973. Outline of a New Approach to the Analysis of Complex Systems and Decision Processes, IEEE Trans. on Systems, Man, Cybernetics, SMC-3(1), 28-44.
17. Procyk, T.J., Mandami, E.H., 1979. A Linguistic Self-Organizing Process Controller, Automatica 15, 15-30.
18. Sugeno, M., Ed., 1985. Industry Applications of Fuzzy Control, North-Holland.
19. Heller, Martin, 1998. At Your Command, Windows Magazine, July.
20. Bakshi, B.R., 1998. Multi-scale Principal Component Analysis with Application to Multivariable Statistical Monitoring, AIChE Journal, July.
21. MacGregor, J.F., Kourti, T., 1998, Multivariable Statistical Treatment of Historical Data for Productivity and Quality Improvements, FOCAP0 98, Snowbird.
22. Hwang, D.H., Ahn, T.J., 1998. Process Operation Improvements Based on Multivariable Statistical Analysis, AIChE (to be submitted).
23. Dudzic, M., 1998. The Use of Advanced Multivariable Statistical Technologies (Chemometrics) at Dofasco, AISE Conference, MIT, MA, July
24. Amari, S, Kasabov, 1997, Brain-Like Computing and Intelligent Information Systems, Springer-Verlag. Singapore.
25. Van Altrock, C., 1997, Fuzzy Logic and Neurofuzzy Applications in Business and Finance, Prentice Hall PTR, New Jersey

TEAM SCHEDULING BY GENETIC SEARCH

Tiehua Zhang*, William A. Gruver*, and Michael H. Smith**

*Intelligent Robotics and Manufacturing System Laboratory
, Simon Fraser University
Burnaby, BC, V5A1S6 Canada

Email: tzhang@cs.sfu.ca, gruver@cs.sfu.ca Web site: <http://www.ensc.sfu.ca/irms>

**University of California at Berkeley
Department of Electrical Engineering and Computer Science
Berkeley, CA. USA
mhs@robotics.eecs.berkeley.edu

ABSTRACT

We consider a photographic studio that must schedule multiple teams of photographers to a large number of elementary and secondary schools. The photographers' schedules are to be optimized so that time constraints are satisfied and each team is able to at least visit two schools daily. A multiple Travelling Salesman model is used where the total distance traveled and time consumed can be evaluated in a single cost function to achieve overall optimality. A genetic algorithm has been applied to solve the problem. The results show that this approach rapidly provides an effective means for solving the problem.

INTRODUCTION

Scheduling involves allocating human operators and material resources to machines in a manufacturing environment so as to meet specified priorities [7] [9]. It attempts to assign and sequence these shared resources so that industrial constraints are satisfied and production costs are minimized.

Effective scheduling can improve on-time delivery of products, reduce inventory, reduce lead-time, and improve the utilization of bottleneck resources. However, in actual industrial practice, more emphasis has been placed on job-shop scheduling and production planning where the task is to allocate machine time and determine a sequence for a set of jobs, each comprising certain operation steps, to be processed on available machines. The objective is either to meet due dates or to minimize the makespan. Production scheduling, however, constitutes only one of several factors in manufacturing which need to be optimized in order to increase production efficiency. In global markets, where competition is fierce and labor costs are high, exploring the potential of people and making staff interact more efficiently in the work environment has become a major issue. This leads to an equally important topic, staff scheduling.

This paper proposes an employee scheduling method using a genetic algorithm as a search method. The formulation of the problem is based on the Travelling Salesman Problem (TSP) model [4]. Our task is to develop a daily schedule of photographic teams to take pictures of school children. The main concern is to satisfy the time constraints so that each day every team can finish at least two schools. We show how to combine a genetic algorithm with a multiple-TSP model.

STAFF SCHEDULING

Staff scheduling, also known as workforce allocation, creates schedules that provide the best possible work coverage and meet employee preferences. The concerns are centered on employees, rather than machines. For example, Mason, *et al* [8] describes a simulation and optimization approach for personnel scheduling of customs staff at the Auckland International Airport. He used a simulation system embedded within a heuristic search and linear programming techniques to determine minimum staffing levels, and then created rosters to ensure the passenger processing targets were satisfied. Weil and Heus [10] applied constraint-based methods to hospital nurse scheduling with a small number of shifts.

In staff scheduling, previous studies such as these treated situations in which the work occurred on-site at company facilities. There are frequent situations, however, when the service has to be offered at clients' facilities and involves travel. Several industries bear similarities with this model. In airline crew scheduling, for example, flight attendants are assigned to the airline's daily service routes. There are location constraints (an employee in Paris cannot board a plane starting from New York to Sydney), as well as federal regulations (an attendant is not allowed to fly more than a certain number of hours per week). Another example is a company that sends its teams to promote new products, or a courier company provides service in a city with its vehicle fleet. In the latter case, the objective may be to minimize the total cost or time to finish the task.

The approach discussed in this paper developed from the authors' collaboration with a photographic studio. Technically, the studio is a manufacturer of high quality photographic products for its customers. In this sense, the studio is like other manufacturers of goods and services. They need raw materials (camera equipment, photo print paper and chemicals) to order from outside sources and to allocate them in the studio. They train photographers, lab engineers and front desk receptionists and schedule their labor force up to 3 shifts per day in the summer season. They provide digitized photo packages in CD-ROM format for special groups such as hockey clubs and schools. Also, they have agreements with public transit authorities to make bus passes for school children. Furthermore, after all pictures are taken, there are many processing operations: rolls of films to be developed in the photo lab, proof prints to be produced, and package ordering sent for.

For years, the studio has been unable to adequately cope with photographic team scheduling. For example, during the Fall semester, eight teams may be sent to 800-1000 elementary and secondary schools (some in remote locations) to take pictures of all children in the schools. These schools, with a student population of 700 to 2300, are booked by customer service early in the year. Each team consists of at least four photographers, four cashiers and a person who coordinates the teachers and students in the school after the team arrives. Each team can usually finish two schools in a day. The team will have to perform equipment set-up, photography, and travel to the two schools. Considering the number of students and time limitation, this is not an easy task. The studio's main concern is to minimize total time. To model this process, we define the following objective to be minimized:

$$F = \sum_1^4 (w_1 s_i + w_2 d_{i1} + w_3 d_{i2}) \quad 1.$$

where d_{ij} is the distance from the studio to a school, s_i is the distance between two schools, and w_i is the weight associated with the travel time. Overall distance is not an issue, but the total time consumed each day is important. If a team is delayed by traffic, other schools will be affected.

In the past, scheduling for dispatching the photography teams has been done manually based on salespersons' experience. Because of the nature of this problem, we choose to formulate it using a travelling salesman problem model explained in the next section.

TRAVELLING SALESMAN PROBLEM

The Travelling Salesman Problem (TSP) is a classical combinatorial optimization problem in which the optimal solution can be determined by exhaustive search. Since it is the mathematical abstraction of many practical application situations, it has been widely investigated in the literature and benchmarked for search algorithm effectiveness. For example, starting from a home point, suppose a salesman wants to visit n cities (n nodes) exactly one time and then return to his home city as illustrated in Fig. 1:

The goal for this TSP is to find the shortest tour sequence so that the distance traveled is minimized. To do this, Noschang [4] described the use of Tabu search, simulated annealing, neural networks, limited memory heuristic search, and genetic and evolutionary algorithms. Louis and Li [3] applied a modified GA to solve the TSP by evaluating the performance of the genetic algorithm with stored data from previously solved similar problems. However, it should be noted that solving a TSP is extremely difficult in many practical situations. For example, the travel time from A to B is not necessarily shorter than A to F, e.g., there may be a bridge between A and B, or heavy traffic. Furthermore, if there are 20 schools to schedule, we obtain

1.22×10^{17} different potential paths. To illustrate this, Table 1 shows a path of a simple TSP of 4 schools, where the numbers indicate distances between schools. It is an asymmetric TSP where, for each team, the travel time from the studio to a school, to the second school, and then back to the studio is dependent on time and geographic locations. It is a multi-salesman problem where several teams share the daily workload and they are expected to achieve a global optimal schedule.

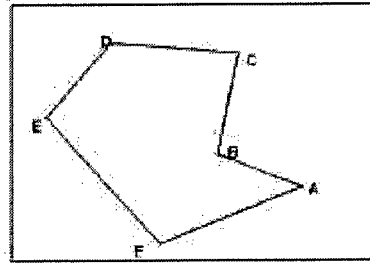


Fig. 1. Travelling Salesman Problem

Table 1: TSP Formulation

(km)	School A	School B	School C	School D	Studio
School A	--	41	16	48	34
School B	41	--	57	46	63
School C	16	57	--	29	52
School D	48	46	29	--	38
Studio	34	63	52	38	--

GENETIC ALGORITHM SEARCH

Since the late 1980s, genetic algorithms [1] have been studied as a search method and applied successfully to a wide range of application situations, including production scheduling and constraint satisfaction [2] [9]. Genetic algorithms are based on the theory of natural evolution where only the most suited individual in a population survives. It manipulates a set of randomly generated potential solutions to the problem. GA provides a powerful search method for optimization and is extremely helpful for problems when little empirical knowledge on the system is available to formulate analytical expressions for the solution.

Computationally, a simple genetic algorithm typically involves four stages:

(1) Binary string encoding. The optimization variables are represented by an encoded string of bits, mappings from each possible solution to a unique binary value. For either continuous values or binary quantities of the potential solutions to a search problem, the encoding scheme uses their integer representation by linearly mapping the variables to an integer in a specified range, and then encoding them using a fixed number of bits. The binary codes of all the variables are concatenated to form a computer representation of the potential solution to the problem. If 16 bits represent a solution for the GA search,

1101011101110100 is a possible string in which the first 8 bits are allocated for x and the last 8 bits are designated for the second variable. For a combinatorial problem like the TSP, a binary string for a solution does not have an obvious meaning. We simply use the natural representation, such as EAHDFGBC for an n -school situation. The first two letters in this tour imply that a team has to start from the studio and first visit school E, then go to school A before returning to the studio. To calculate the cost, time weights and distances from the studio to each school are added.

(2) Solution evaluation. The objective function to be optimized gives a means for evaluating each string in the population. The fitness value is often normalized to maintain uniformity over various problem domains. As there is an obvious meaning to the TSP evaluation function, we consider it as a floating-point number.

(3) Selection. Production of later offspring is based on their fitness values in the population. A fitter solution/string receives a higher chance of surviving in the subsequent generation. Usually a selection scheme allocates offspring based on the ratio of a string's fitness value to the population's average value. Statistically the allocated number of offspring approaches the expected number only for very large population sizes.

(4) Genetic manipulations. Crossover is the process of picking pairs of candidates from the population to exchange information at a randomly chosen point between them. This is controlled by the crossover rate activated if a random number is greater than it, otherwise the strings are transferred to the next generation. After crossover, mutation is performed by bit flipping changes in a string, which is controlled by the mutation rate. This attempts to restore lost genetic material in the process of evolution. These two operations cannot be directly applied to a TSP. For example, if the two following parents are chosen by crossover operation, and the exchange point falls at the fifth position:

EAHDF|GBC
DBFHA|CEG

This results in two children of:

EAHDF|CEG
DBFHA|GBC

which are both invalid solution representations because school E appears twice in the first child and B is never visited. The second child is not valid either. For mutation, there is no flipping to an alphabetical string. None of the schools in the above encoding scheme can be flipped to restore the possibly lost genetic information. These two aspects will be discussed in the following section.

Finally, termination of the program depends on choice of a stopping criterion. Termination could occur after a fixed number of generations, after a string with a certain fitness value is obtained, or after all the strings in the population have attained a certain degree of homogeneity. Here, as in any optimization process, we have to consider the possibility of premature convergence [3].

IMPLEMENTATION

To illustrate our approach, we use a simple genetic algorithm to implement a search for team scheduling for 4 teams. Since the TSP has a special combinatorial structure, there are special features for its GA operations:

Since binary chromosomes cannot be used for string encoding, we employ a direct character array representation, such as DCAFBE, which depicts a tour of six schools.

For TSP crossover operation, we used the Grefenstette Greedy Crossover [3][5]. Greedy Crossover has a four-stage process:

1. For each pair of parents, pick a random school to start;

2. Compare the two routes leaving the school and choose the shorter route;
3. If the shorter parental route would introduce cycles into the partial tour, then extend the tour by a random route;
4. Continue to extend the partial tour using steps 2 and 3 until the tour is complete.

The heuristic is "greedy" since it always selects the locally shorter route, which is based on specific knowledge of the TSP. Locally, shorter paths are better than longer paths in most cases. In the four-team scheduling situation modeled as a multiple TSP problem, there are four subtours to search in order to achieve a globally optimized solution. Nevertheless, it does not necessarily lead to a more complex problem.

In the crossover process, suppose we have two parents:

EAHDFGBC
DBFHACEG

To obtain an offspring, we select school E as the starting point for the first child:

E - - - - -

Then we compare the distances of the two schools in the parents leaving from school E, which are EA and EG, and choose the closer school as the second one to visit. If EG is shorter, we obtain

EG - - - - -

At this point, we have two options of the first school for the second team, H and F. We randomly select F and obtain:

EGF - - - - -

Next, we compare the edges leaving F, FH and FG. Since G is already in this tour, choose H,

EGFH - - - -

Continuing this procedure, we complete all the schools for four teams and assign it as the first child. The second child can be obtained using the same method.

For TSP mutation, we choose two schools at random from a parent and then swap these two schools. If the second and the fourth school are chosen for mutation in the parent, we obtain a child as shown below. We maintain a mutation rate of 0.067 during the search process:

EAHDFGBC -- prior to mutation

EDHAFGBC -- after mutation

There are other reported selection heuristics, such as "Keep the Best" (KTB), where individuals from both parent generation and child generation are evaluated and the most suitable ones in the whole group are retained. For simplicity, roulette wheel selection is applied. Our program typically can find a reasonable solution in a timely fashion. Time dependency is not considered at this stage, but the next step is to implement a hybrid GA search algorithm which addresses the issue of optimizing different types of cost functions.

CONCLUSIONS

The paper presents an approach for solving a multiple team-scheduling problem using a Travelling Salesman model and applying genetic search techniques. The significance of using a genetic algorithm in calculating a solution to this n-school scheduling problem may not be obvious. Indeed, many AI techniques

could also solve this problem. However, when the number of teams and schools increases and the task becomes more complex, the search time taken by a GA will not exponentially increase. Moreover, after a workforce allocation can be modeled, the proposed method can be easily modified to perform the search.

The main goal in this phase of the project is to replace the manually generated non-optimal schedule. Other goals include investigating methods to integrate employee experience, qualifications, availabilities, and pay-rates into a system. Since there is always the possibility of machine breakdown or a team member may become sick on a particular day, it would be important to have a means for rescheduling the team without changing the contract with a school. Also, scheduling methods such as these must be integrated with other parts of the manufacturing system, e.g., enterprise resource planning, machine controls, and maintenance, in order to form an information infrastructure that will improve the overall efficiency of the studio. Finally, the approach that is employed here to solve this class of problems applies to other similar tasks and manufacturing processes. For example, multiple robots can be scheduled to perform welding or, in circuit board manufacturing, an optimal sequence for component placement can be planned.

ACKNOWLEDGEMENT

This research received support from the National Science and Engineering Council of Canada. The authors thank David Robinson, Amber Computer Systems, Inc., Delta, BC, for his collaboration on this research.

REFERENCES

1. D.E. Goldberg, 1989. *Genetic Algorithms in Search, Optimization, and Machine Learning*, Addison-Wesley Publishing.
2. D.C. Mattfeld, 1996. *Evolutionary Search and the Job Shop; Investigations on Genetic Algorithms for Production Scheduling*, Springer/Physica Verlag.
3. S. J. Louis, Gong Li, 1997. "Augmenting Genetic Algorithms with Memory to Solve Travelling Salesman", Proc. of the Joint Conference on Information Sciences, Durham, NC.
4. M. H. Noschang, 1997. "The Travelling Salesman Problem – A Review of Theory and Current Research", College of Engineering, University of Cincinnati, Technical Report.
5. V.M. Kureichick, V.V. Miagkikh, A.P. Tochy, 1996. Genetic Algorithm for Solution of the Travelling Salesman Problem with New Features against Premature Convergence", Proc. ICGA96, Gursuf, Ukraine.
6. J. Grefenstette, R. Gopal, R. Rosmaita, D. Gucht, 1985. "Genetic Algorithms for Travelling Salesman Problem", Proc. of the Second International Conference on Genetic Algorithms, Mahwah, NJ.
7. F.A. Rodammer, K.P. White, 1988. "A Recent Survey of Production Scheduling", IEEE Transactions on System, Man and Cybernetics, 18(6).
8. A.J. Mason, D.M. Ryan, D.M. Panton, 1988. "Integrated Simulation, Heuristic and Optimization Approaches to Staff Scheduling", Operations Research, 46(2).
9. F. Archetti, M. Lucertini, P. Serafini, 1989. *Operations Research Models in Flexible Manufacturing Systems*, Springer.
10. G. Weil, K. Heus, P. Francois, M. Poujade, 1995. "Constraint Programming for Nurse Scheduling", IEEE Engineering in Medicine and Biology, 2.

Intelligence in Surface Processing of Materials

An Intelligent AE Sensor for the Monitoring of Finish Machining Process

S. Dolinšek*, J. Kopac*, Z.J. Viharos, L. Monostori ****

*University of Ljubljana, Faculty of Mechanical Engineering, Ljubljana, Slovenia

** Computer and Automation Research Institute, Hungarian Academy of Science,
Budapest, Hungary

ABSTRACT

The paper presents the latest results of sensing the cutting process on the basis of AE signals and some particularities in further development of the monitoring model for the finish turning process. Due to non-linearity, the large number of influencing parameters and missing information in AE data, the Artificial Neural Networks were chosen as a monitoring decision tool. The problem of accurateness in predicting the surface roughness on the basis of AE - because of the mutual interdependence of the data - requires a special procedure for building a neural network model. The final aim of such an approach is presented as improvements in learning or considerable reduction in error prediction. Further development of the monitoring model has the goal of building a so-called intelligent sensor, which should be able to perform the signal conditioning and feature extraction process.

INTRODUCTION

Most of the reports on research into the machining processes usually start with a similar ascertainment: the complexity of the cutting process is one of the main obstacles to successful modeling or monitoring of processes; this fact gives us the impetus to continue permanent investigations. There are no simple answers or quick solutions; a reliable monitoring approach or a complete control system for the cutting process is a task in which successful solutions could be obtained only through numerous, systematic investigations covering the different scientific areas incorporating sensor technologies, signal processing techniques, modeling methods, etc.

Probably one of the best of the latest reviews of such efforts has been made within the CIRP groups, where the conclusions stated that the different monitoring systems with acceptable commercial reliability are now available in the market, although the narrow range of performance provides only limited applicability (Byrne and others at [1]). The report also confirms one of the main gaps in this kind of research; i.e. to develop a system as an integrated part of an intelligent machine tool, much more should be done at both the hardware and software level to obtain a simple and reliable sensor for machining applications. At present - when the development of manufacturing processes heavily depends on information technologies - the realistic process models are also one of the prerequisites for predicting the performance of metal cutting operations. The latest report on the modeling of machining operations (C.A. van Luttervelt and others at [2]) concluded, that most of the research deals with possible new ways of obtaining better control of machining operations; however a common framework is still missing.

At the first IPMM conference the monitoring concept on the basis of sensing Acoustic Emission signals in finish machining processes was presented (see Dolinsek at [3]). From the contents of the AE signals we were able to extract significant features from the process, depending on the cutting conditions, which serve as learning data for the ANN structure. The model should be applicable for practical cutting in such a way that the predicted values of surface roughness could be a sign to adapt the cutting parameters in order to achieve the required surface quality or to detect disturbances in the process (tool wear, unfavorable chip shape, lack of coolant). In the introduction we also draw attention to the lack of adequate sensors and indicate that the sensing technology will play an important role in the development of future manufacturing systems.

Further investigation of our monitoring concept for finish machining processes was therefore oriented towards the search for reliable sensing. Some of the results using the AE-jet sensor were discussed at SEM and CIRP conferences (see Dolinsek at [4] and [5]). The main advantages of this sensor were presented as improvements of the signal to noise ratio, simple upgrade, and the fact that the cutting process and sensor are not reciprocally disturbed. Through the spectral analysis technique, and with adequate averaging procedures, we were therefore able to gain some useful information for the further development of our monitoring model.

Artificial neural networks (ANNs) were used as a operating tool because they can handle strong non-linearities, a large number of parameters, missing information, and the characteristics of the data which are also significant in our monitoring approach. Based on their inherent learning capabilities, ANNs can adapt themselves to changes in the production environment, and can also be used in case where no detailed information is available about the relationships among the various manufacturing parameters. In many cases, there is also no exact knowledge about the relationships among parameters; it is unknown which input-output configuration of an ANN can satisfy the accuracy requirements. Therefore, a method is needed for automatic input-output configuration of the applied ANN model. This paper therefore addresses the problem of automatic input-output configuration and generation of ANN-based monitoring models, i.e. those parameters to be considered as inputs, and those as output, in order to accurately predict surface roughness and classify tool wear in the finish turning process.

AE-JET SENSOR FOR MONITORING A FINISH MACHINING PROCESS

In researching Tool Condition Monitoring (TCM) systems for the manufacturing processes and introducing them to the workshop environment we are engaged in solving three main tasks :

- building up a sensor system which is reliable for sensing the process parameters with minimal influence on the process,
- applying proper signal processing techniques capable of processing the real life signals,
- developing decision-making algorithms capable of estimating the process conditions.

In such an feature-based approach, we observe some features, extracted from sensor signals in order to identify different process conditions and compare them to normal and unfavorable cutting conditions. This process is generally not too complicated, but the success of the monitoring depends greatly on exact correlation of the measured parameters to the cutting process characteristics - i.e. the sensors are the first and main component leading to the successful solution of our tasks. Once we find or build-up a sensor which satisfies the main requirements demanded in practical monitoring approaches: measurement close to the machining point; no influences on the machine-tool characteristics; function independent of tool or workpiece; low costs; maintenance and wear free; resistance to dirt and to mechanical and thermal influences, minimal reciprocal disturbances between the process and sensor, simple upgrade which allows easy further improvements. Thus we can further develop our monitoring model by applying signal processing, feature extraction and decision making procedures. When this intelligent part is successfully solved the final hardware integration into the sensor is not a complicated task.

One of the most promising tool monitoring techniques is based on sensing the Acoustic Emission (AE) signals generated at the cutting zone. Extensive publications have demonstrated the extreme sensitivity of AE signals to certain process parameters (etc. Dornfeld at [6]). In general it is agreed that during metal cutting, plastic deformation (continuous type of AE signals) and fracture of the material (burst type of AE signals) are major sources for AE waves. One of the basic researches of the AE phenomena in the cutting is that made by Moriwaki, who illustrated in detail seven possible sources of the AE signals in the vicinity of the cutting process [7]. However in sensing and analyzing the AE signals generated in the cutting process we will always face two main obstacles:

- it is almost impossible to built-up a physical model of the AE signal in relation to the AE waves, since the signal generated in real cutting processes in complex workpiece structures is continuous and random,
- we expect the sensor used in metal cutting problems to be reliable in sensing AE signals from all sources of generation with the ability to differentiate between the sources, but without any other interferences.

Due to those limitations, the most common approach of the application of the AE in monitoring of cutting processes is at present still a simplified sensing of the mixed AE signals. From the content of the acquired AE signals using a suitable post-processing procedure one can then identify different process conditions. Although many different sensors are available for AE measurements, only few can be used in a machine tool in which aggressive ambient conditions occur. The main disadvantages of traducers, which are mainly designed for non-destructive inspection or research work, are that they cannot withstand the high temperatures, large coolant volumes and abrasive wear through chips. With a new concept of AE transducers (see [8]), - the Water-jet AE sensors - a liquid or coolant stream is used as a transmission medium to transfer AE signals generated from the cutting process to the PZT element. As the distance between the cutting zone and transducer element is small, the damping effect is minimized, and the signal-to-noise ratio is significantly improved. The construction of the sensor, developed for our monitoring task, presented in Fig 1, was a practically built-in CNC finish turning machine. The applicability of this sensor in the finish turning process was tested throughout the proper analysis of the acquired AE signals and further relation of their content to the process conditions.

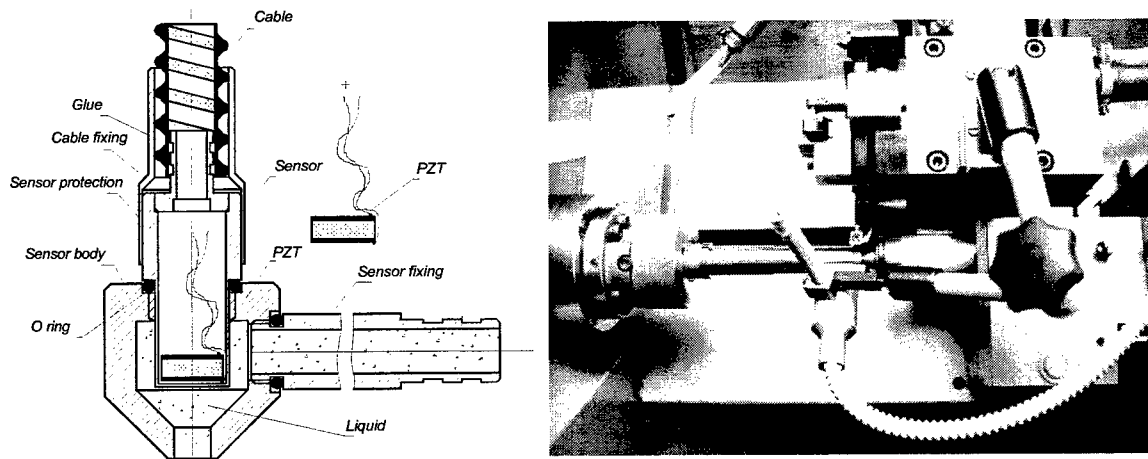


Fig. 1. AE-jet sensor for monitoring a finish machining process built-in CNC turning machine

All the tests were made by using cutting conditions producing continuous chips. Therefore, a frequency analysis of the signals - using the well known assumption for continuous random signals - could be applied as the signal processing technique. In order to obtain satisfactory amplitude estimations in spectral analysis, 75 averages within sample and 15 averages between the samples, were performed. Fig. 2 presents the response of the sensor when the tool is not cutting (free run of the turning machine). We can see that the power spectra of the signals shows a distinctive amplitude peak in the range of the resonant frequency of the sensor, and that the energy of the signal is mainly distributed within the range of 100-610 kHz.

In Fig. 2, which also presents a comparison of the spectra of AE signals obtained in cutting with different parameters, we notice that the spectra and their energy are altered according to the changes in cutting conditions. From the results presented here, and those already published [4,9], the following conclusions can be drawn:

- the energy of the AE signal is mainly distributed in the frequency range of 100-610 kHz,
- the sensitivity of the sensor depends on the position of the piezo-ceramics placed inside the coolant stream,
- this sensor can be used in workshop conditions with the coolant supplied by the machine tool pump,
- tool wear is one of the most influential factors in increasing the energy of the AE signal.
- the sensor is sensitive to AE signals obtained by cutting different workpiece materials and according to the variations of the machining parameters; their influence can be seen from different spectral energies.

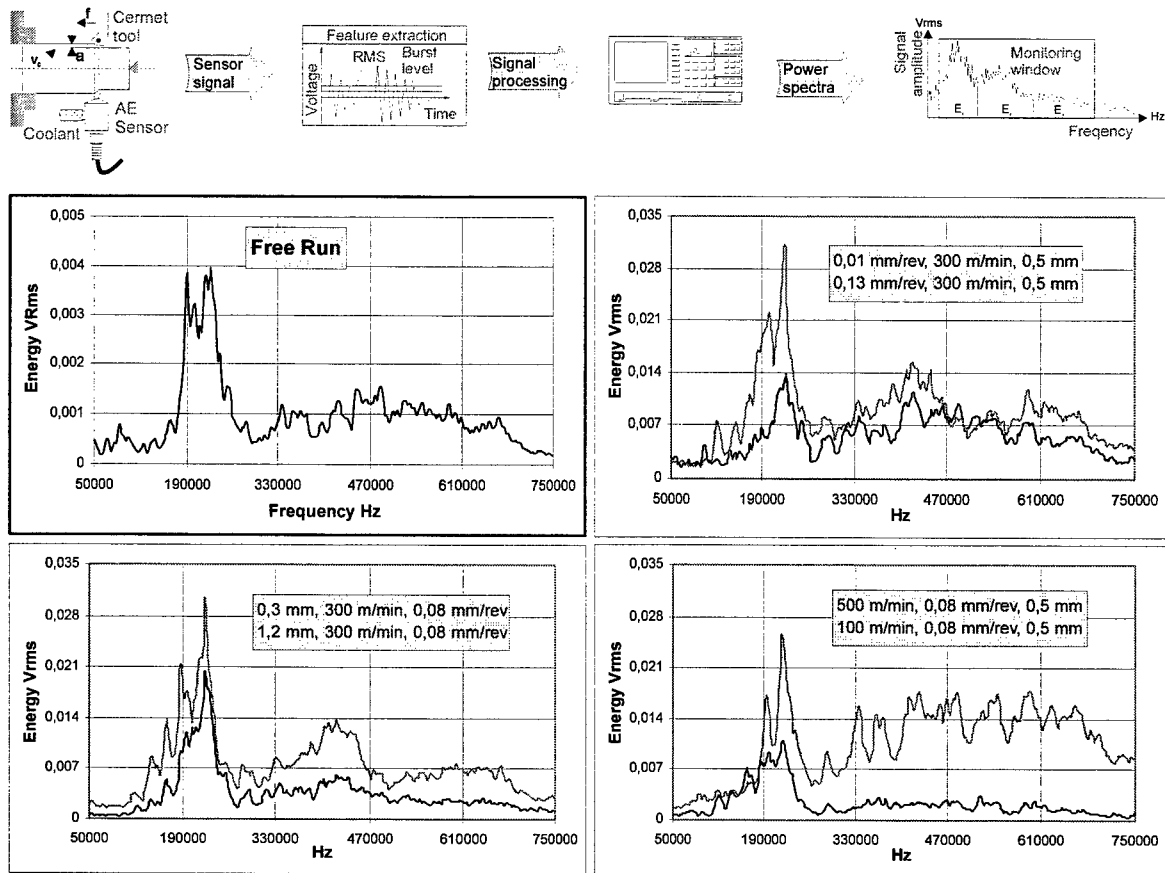


Fig. 2. Spectra of the AE signals in finish turning process obtained in cutting with different cutting conditions

ANN AS A TOOL FOR INTELLIGENT PROCESSING OF AE SIGNALS

Considering the input and output variables in our model for monitoring a finish turning process as a set of parameters, the ANN model estimates a part of this parameter set based on the remaining part. This selection strongly influences the accuracy of the developed model, especially if dependencies between parameters are non-invertible. In different cutting conditions (e.g. first cut with sharp tool, non-uniformities in material structure, disturbances in coolant flow), the tasks are different; consequently, the estimation capabilities of the related applied models are different, even if the same set of parameters obtained with the same cutting parameters is used. One of the main goals of the research was to find a general model for a set of assignments which can satisfy the accuracy requirements. This goal was achieved by a sequential forward selection (SFS) search algorithm, which uses the heuristic of speed of ANN learning. This method incorporates:

- determination of the number of output variables,
- determination for each parameter to be either input or output.

The method also builds up the appropriate ANN model without considering the assignment of an engineer; it is also useful in the case of strong non-linear relationships. [10]. Research also focused on how to apply the general model for various tasks. Usually, the engineer knows some parameters of a process and the modelling task is to determine the other parameters while satisfying some constraints. After obtaining the general ANN model, in almost every case, part of the input and output variables of the general model are known by the user; the task of modelling is then, to search for the remaining, unknown input and output parameters. In order to realise this, a simulated-annealing search method was used to determine the unknown input parameters. After obtaining the appropriate input parameters, the unknown output

parameters can be determined by a simple ANN estimation. The values of the unknown input parameters are appropriate, if:

- they are between their minimum and maximum values,
- the estimated unknown output parameters are between their minimum values,
- the estimation of the output parameters determined on the basis of known and unknown input parameters ensures that the estimated values of the known output parameters are equal to their known values

The first and second conditions determine the validity of the ANN model. With the help of this method, all possible assignments of an engineer can be solved by the general ANN model [11]. As a practical demonstration of the method, an example of the evaluation of the AE signals from finish machining experiments will be analysed. In this investigation, the target is to estimate the roughness of a product surface based on known values of cutting parameters and measured AE signals. The cutting parameters, which were varied, are feed ($f = 0,01-0,2$ mm/rev), cutting speed ($v_c = 100-500$ m/min) and depth of cut ($a = 0,1-1,2$ mm). From the measured AE signal and its energy content, significant features were calculated using particular energies found in four frequency ranges (according to the findings from previous investigations [3,4]), 50-750 kHz, 99-249.5 kHz, 249.5-400 kHz and 400-610 kHz. These seven parameters together act as known parameters, while the roughness acts as the parameter to be estimated. Since the non-linear dependencies among these parameters are already experienced, an ANN model can consequently be used to realise the mapping among parameters.

The concept of building an intelligent sensor for a finish machining process is based on the idea of performing on-line learning as fast as possible over the whole range of applicable cutting parameters with satisfactory accuracy in prediction of the process conditions. Using statistical planning of experiments, 21 measurements were used to build up an ANN model (the selection of parameters are presented in Fig. 3) with seven machining situations to test the behaviour of the applied ANN model. The above-described method was used to demonstrate the automatic input-output configuration of the ANN.

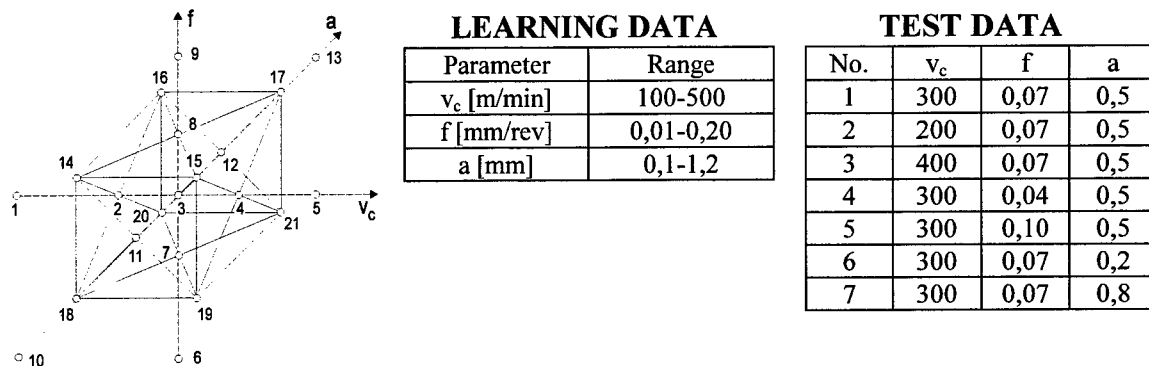


Fig. 3. Planning of the experiment for learning and test data.

Because the number of learning vectors is small, we firstly tried to build up an ANN model with one hidden layer and two hidden nodes. The target average estimation accuracy of the ANN model was $\pm 2.5\%$. The above described method found one output parameter that can be estimated by the ANN model based on the remaining parameters (inputs). This was the energy parameter E_1 (50-750kHz). The roughness becomes the input of the model. Fig. 4 shows the input-output configuration of this ANN model. With the help of this new method it is therefore possible to estimate the unknown parameters based on the values of known parameters, regardless of whether were the input or output of the ANN model. In such a way, the estimation of the roughness based on known values of v , f , a , E_1 , E_2 , E_3 , E_4 was performed. The seven test situations were reviewed at the first stage. In each of the situations, the estimations were repeated ten times to check if there were more solutions for the given estimation task. The estimations of surface roughness with this ANN model, presented in Fig. 4, are not very accurate, but the developed algorithm reports through the parameter E_1 that the roughness can not be estimated accurately enough.

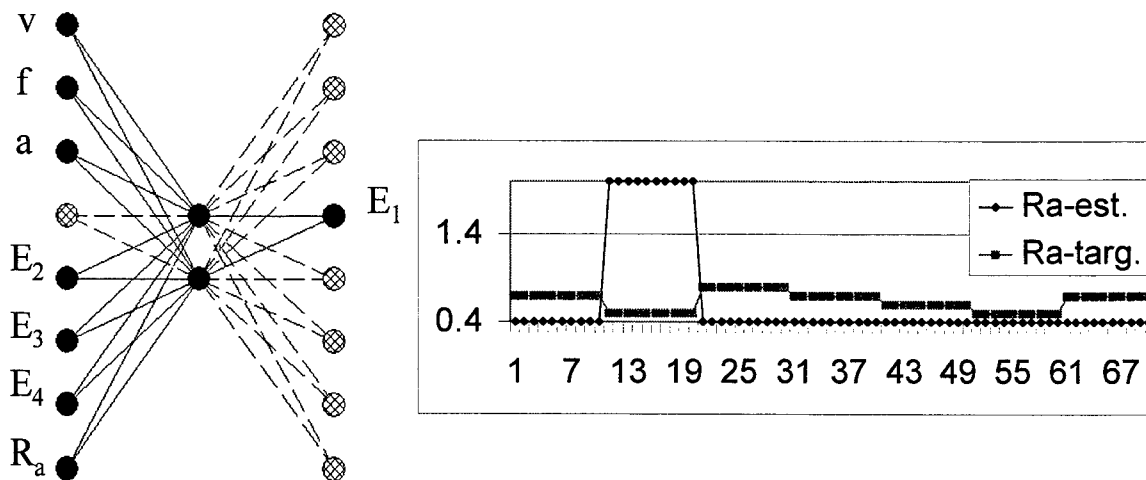


Fig. 4. The input/output configuration of the first ANN model with estimated and target roughness.

Very similar estimates were obtained with the situations used to build the ANN model. The conclusion is that the ANN model is inappropriate; consequently, the number of input nodes must be enlarged. A second investigation was performed with an ANN having 6 hidden nodes. In this case, the I/O configuration of the resulting ANN was different from the previous one. This ANN with five inputs and three outputs is presented in Fig. 5. The outputs are E₁, E₂, and E₄. To check if the model was sufficiently accurate, the roughness estimation was performed on the learning data set. The results of these estimations show that the ANN model learns the dependencies between the inputs and outputs from the learning data set. To test the model in test cutting situations, estimates of roughness were performed with this new ANN configuration, repeated ten times for each situation. The results of the estimated and measured roughness (Fig. 5) show that the ANN is unable to estimate the roughness in these situations, although it could perfectly estimate the learning data. This shows that in further investigation the quantity of learning data must be increased.

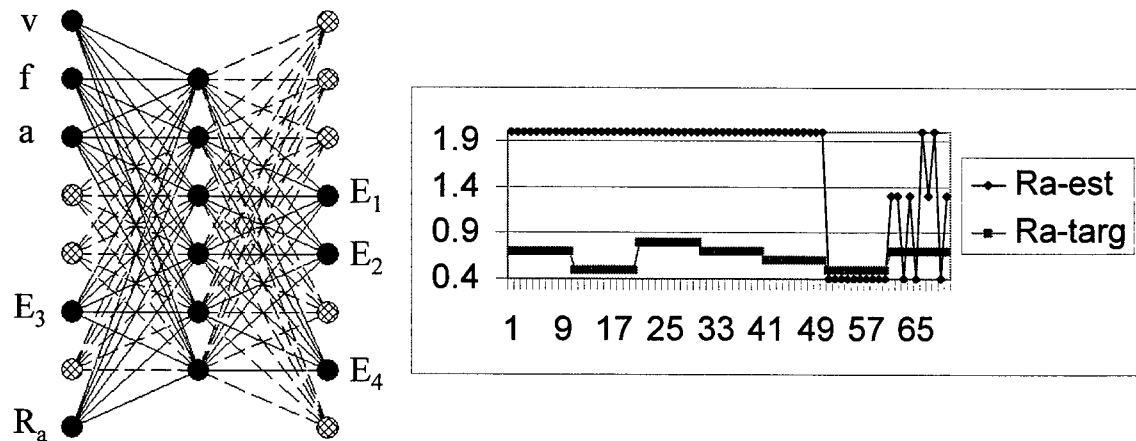


Fig. 5. The input/output configuration of the second ANN model with estimated and target roughness.

CONCLUSIONS

The introduction of automatic control of the machine tools, and the highly invisible cutting area inside the machine tool, have made monitoring of the processes even more difficult. In recent years, extensive efforts have been spent to develop such systems; however, reliable and marketable products are still unavailable. In this paper, the latest results are presented in using a more sensitive sensor - the AE-jet sensor - to monitor the finish turning process in order to obtain the prescribed surface finish of the products. The response of the sensor is adequate to the damping characteristics of the cutting process in relation to changes in cutting

conditions; response can be observed directly in spectral energies at different frequency ranges of the AE signal. Using the energies of the AE signals, obtained during cutting, under both favorable and unfavorable conditions, we have built-up the concept of monitoring decision-making. Artificial Neural Networks can handle non-linearities, multidimensionality, missing information and also the main problem in the AE results obtained, the overlap presented by the same values of roughness or AE for different cutting parameters. With the new method of building an ANN structure, which can determine the number of variables and search for input and output variables, the learning phase is shortened and the prediction is improved. The results obtained demonstrated a good direction for further work in this direction.

REFERENCES

1. G. Byrne, D. Dornfeld, I. Inasaki, G. Ketteler, W. König, R. Teti, 1995. Tool Condition Monitoring, The Status of Research and Industrial Application. *Annals of the CIRP*, 44(2), 24-41.
2. C.A. van Luttervelt, T.H.C. Childs, I.S. Jawahir, F. Klocke, P.K. Venunod, 1999. Present Situation and Future Trends in Modeling of Machining Operations. *Annals of the CIRP*, 48(2), 587-626
3. S. Dolinsek, 1997. Hybrid intelligent systems in monitoring of finishing machining processes, Australasia Pacific Forum on Intelligent Processing and Manufacturing of Materials. Brisbane, Australia.
4. S.Dolinsek, 1997. Tool Condition Monitoring Using AE-jet sensor, 1997. Postconference Proceedings of the SEM Spring Conference, SEM , 165-174.
5. S. Dolinsek, 1998. Monitoring of finish machining process using AE-jet sensor, 31st CIRP International Seminar on Manufacturing Systems, Berkeley, USA, 540-545.
6. J.Liu, D. Dornfeld, 1996. ASME, Journal of Engineering for Industry, 118 (4), 199-207.
7. T. Moriwaki, 1983. Application of Acoustic Emission Measurement to Sensing of Wear and Breakage of Cutting Tool. *Bull. Japan Soc. Prec. Eng.*, 17 (3), 154-160.
8. I. Grabec, W. Sachse, 1991. Automatic Modeling of physical phenomena, Application of ultrasonic data. *Journal of Applied Physics*, 69 (9), 24-32.
9. S.Dolinsek, J. Kopac, 1999. Acoustic emission signals for tool wear identification, 12th International Conference on Wear of Materials, Atlanta, USA
10. Zs. J. Viharos, L. Monostori, S. Markos, 1999. Selection of input and output variables of ANN based modeling of cutting processes. CIRP Proceedings of the X. Workshop on Supervising and Diagnostics of Machining Systems, Poland.
11. Zs. J. Viharos, L. Monostori, 1999. Automatic input-output configuration and generation of ANN-based process models and its application in machining. Proceedings of the XII. International Conference on IEA/AIE Systems, Cairo, Egypt.

A New Fuzzy-Fractal Approach for Surface Quality Control in Intelligent Manufacturing Of Materials

Patricia Melin and Oscar Castillo

Computer Science Department, Tijuana Institute of Technology
P.O. Box 4207, Chula Vista CA 91909, USA
Email: ocastillo@mail.tij.cetys.mx emelin@mail.tij.cetys.mx

ABSTRACT

In this paper we describe a general method to automate quality control in the manufacturing of materials using a new fuzzy-fractal approach. We also show how to implement this new method in an intelligent system to achieve automated quality control in practice. Engineers deal with surfaces and surface properties in a wide variety of contexts. In manufacturing, the goal is to produce a surface with specific physical or chemical properties. The concept of the fractal dimension can be used to classify a complex geometrical object [1]. In this case, we use the fractal dimension to characterize surface roughness of materials for manufacturing applications. On the other hand, we used Fuzzy Logic [12] techniques to simulate the expert evaluation/decision process to obtain the quality of manufactured materials. Quality evaluation is simulated in an intelligent system using as input the information about material roughness and porosity (fractal dimension), and then by applying a set of fuzzy rules to decide, on the degree of quality of the production.

INTRODUCTION

In this paper we describe a new method for surface quality control in intelligent manufacturing of materials based on a new fuzzy-fractal approach. Recently, considerable progress has been made in understanding surfaces through application of fractal concepts. The fact that surfaces are fractal was pointed out by Mandelbrot [5]. This in turn prompted the development of the dynamic scaling approach for describing not only the morphology, but also the dynamics of fractal surfaces [11]. A laser scanner microscope can be used to obtain the geometrical information from the samples of materials extracted from production lines.

We can use the concept of a fractal dimension to classify surfaces according to their geometrical complexity. We can define a set of linguistics for surface roughness and porosity, and then a fuzzy rule base relating surface geometry to corresponding quality values. The fuzzy-fractal approach can be implemented as an intelligent computer program to automate the quality control in material processing.

Fuzzy Logic and Fractal Theory can increase the efficiency (in accuracy and time) of quality control, because an intelligent system has the mathematical algorithms (for fractal dimension) needed to classify the roughness of the material, and also because the intelligent system has the knowledge to decide on the final quality of the manufactured material. In this paper the authors have successfully generalized their previous work on this matter [2, 3], by using the fractal dimension to perform automated roughness classification and by developing a knowledge base for evaluation of production quality using Fuzzy Logic techniques.

FRactal CHARACTERIZATION OF SURFACES

Engineers deal with surfaces and surface properties in a wide variety of contexts. In some applications the goal is to produce a surface with a specific physical or chemical property, but often surfaces are inherently formed in industrial and natural processes. Due to the generality and importance of these processes, developing an efficient approach to characterize surface structure and its dynamics and understanding how surfaces are formed is a challenging problem of practical interest to engineers. Recently, considerable progress has been made to understand surfaces through application of fractal concepts [5] and dynamic scaling theory [11]. For example, a tungsten oxide surface layer exhibits geometrical scaled properties suggesting a fractal structure over a scale range [6]. The fractal dimension of the surface is defined as:

$$d = [\ln N(r)] / [\ln(1/r)]$$

where $N(r)$ is the number of boxes covering the surfaces and r is the size of the box. An approximation to the fractal dimension can be obtained by counting the number of boxes covering the surfaces for different r sizes and then performing a logarithmic regression to obtain d (box counting algorithm).

The fractal dimension can be used to characterize surface roughness of materials. The reason for this is that the fractal dimension measures the geometrical complexity of objects. In this case, surface roughness can be classified by using the numeric value of the fractal dimension. If we consider surface roughness as a linguistic variable and we assign it the following linguistic values: bad, regular and, good we can design a classification scheme as shown in Table 1.

Table 1. Fuzzy rule base for surface roughness.

IF	THEN
fractal dimension low	surface roughness good
fractal dimension medium	surface roughness regular
fractal dimension large	surface roughness bad

The reasoning behind this classification scheme is that when the surface is smooth the fractal dimension of the surface will be close to one. On the other hand, when the surface is rougher the fractal dimension will be close to a value of two. We can define membership functions for the linguistic values of the fractal dimension and the surface roughness considering the range of numeric values of these variables.

This fuzzy-fractal characterization has the advantage of enabling the management of uncertainty in this domain of application. Also, it can be used for quality control in manufacturing of materials, because surface roughness is one of the most important properties in material processing. We can also apply the concept of fractal dimension to characterize surface porosity in a similar way. By considering porosity as a linguistic variable, a similar classification scheme can be used with a different membership function scale.

FUZZY LOGIC FOR QUALITY CONTROL

Fuzzy logic techniques can be used to achieve automated quality control in material processing. In this case, a set of fuzzy rules is needed to relate relevant physical characteristics of the material (to be produced) with the quality of the product. In our approach, we assume the fractal dimension can be used as a classification technique which can be input to the fuzzy rules for quality control. The fuzzy rules contain the knowledge of human experts for the specific application domain of material processing. If we consider the quality of the product as a linguistic variable with linguistic values: bad, regular, good, and excellent, we can establish a fuzzy rule base for quality control as shown in Table 2 (only part of the knowledge base is shown). We also use the temperature during the material processing as a linguistic variable with values: high, medium, and low, and the total time of the process with values: large, medium and short.

Table 2. Sample fuzzy rule base for quality control.

IF	AND	AND	AND	THEN
surface roughness	porosity	temperature	time	quality
good	good	medium	medium	excellent
good	regular	medium	medium	good
regular	good	medium	medium	good
good	good	high	medium	good
good	good	medium	large	good
good	regular	high	medium	regular
regular	good	medium	high	regular
bad	bad	low	small	bad

The reasoning behind this fuzzy rule base is that if surface roughness is good and porosity is good, then product quality can be considered good. Otherwise, it can be considered regular or bad as shown in the rules. The use of fuzzy logic in manufacturing applications has been well recognized [7, 8, 9] and many applications have been developed. In this case, we arrived to the conclusion that the best way to convey the information about the quality level of the manufactured product is to use fuzzy sets [12]. Also, we think that the best way to reason with uncertainty in this case is by using fuzzy logic. We can say then that the integration of the power of fuzzy logic with the mathematical concepts of fractal theory enables automated quality control in materials processing.

NEW METHOD FOR SURFACE QUALITY CONTROL USING FUZZY-FRACTAL APPROACH

The new method for surface quality control consists in the integration of the method for fractal characterization of surfaces and the use of fuzzy logic techniques for quality evaluation of the product. We show in Figure 1 how the method works, beginning with the samples of materials extracted from production lines and ending with the final evaluation of production quality.

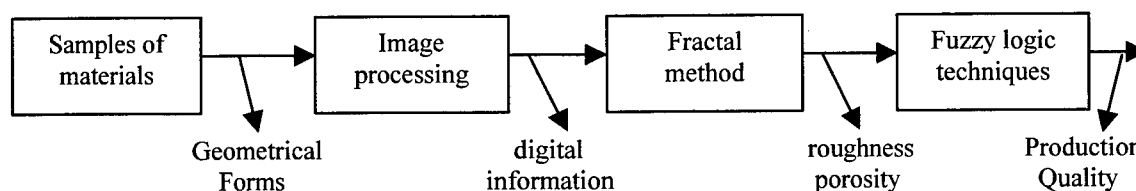


Fig. 1. New method for automated quality control in material processing.

The samples of materials are extracted randomly from production lines (this is the only part that is not automated) and prepared for identification. Then a laser scan microscope is used to digitize the geometry of the surface. The digitized information is then used as input for the method of characterization of the surface (using the fractal dimension). Finally, the fractal characterization of the surface is used as information by a knowledge base of fuzzy rules to decide on the general quality of the production.

We will consider briefly the application of our new fuzzy fractal approach for the case of an oxide layer surface [6]. In Figure 2 we show the oxide layer structure. The oxide depth is divided in two layers. The lower one D2, close to the tungsten substrate C is considered as an homogeneous layer. The upper one D1, limited by the rough surface is a mixed layer filled with a volume V_x of oxide X and a volume V_e of empty space S. The upper layer is the one that is considered for classification purposes using the fractal dimension. The limiting curve (boundary) of the upper layer is used to measure the roughness of the surface as described in Section 2.

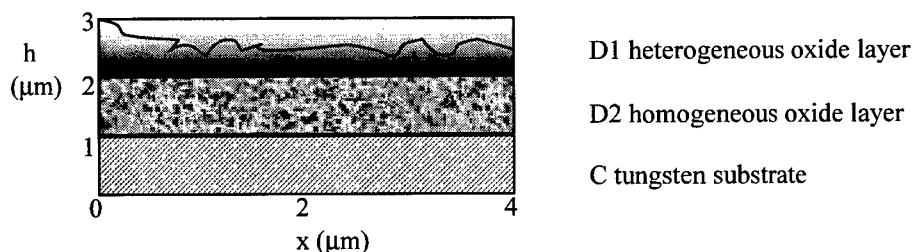


Fig. 2. Oxide profile on tungsten surface.

The new method for quality control simulates the expert decision process involved in obtaining the degree of quality of the production. This method uses as input the information obtained by the fractal method and then applies heuristics of the experts (implemented as fuzzy rules) to decide on the quality of the production.

We have implemented the new method for automated quality control in the MATLAB programming language. The computer program can be considered an intelligent system for quality control in the manufacturing of materials. The two main modules of the intelligent system are the fuzzy logic module and the fractal module. The fractal module consists of a computer program that is an implementation of the method to characterize surfaces using the fractal dimension. This computer program uses geometrical form of surface (obtained from samples of the material) to estimate the fractal dimension (box dimension) using a well known algorithm [5]. The expert module is a computer program that is an implementation of the method to perform automated evaluation of production quality. The knowledge base of this module consists of a set of fuzzy rules containing the knowledge of the human experts for the domain of quality control in the manufacturing of materials.

EXPERIMENTAL RESULTS

To give an idea of the performance of our fuzzy-fractal approach for automated quality control, we show below simulations results obtained for several types of materials. First, we show in Figure 3 the fuzzy rule base for prototype intelligent system developed in the fuzzy logic toolbox of the MATLAB programming language [10].

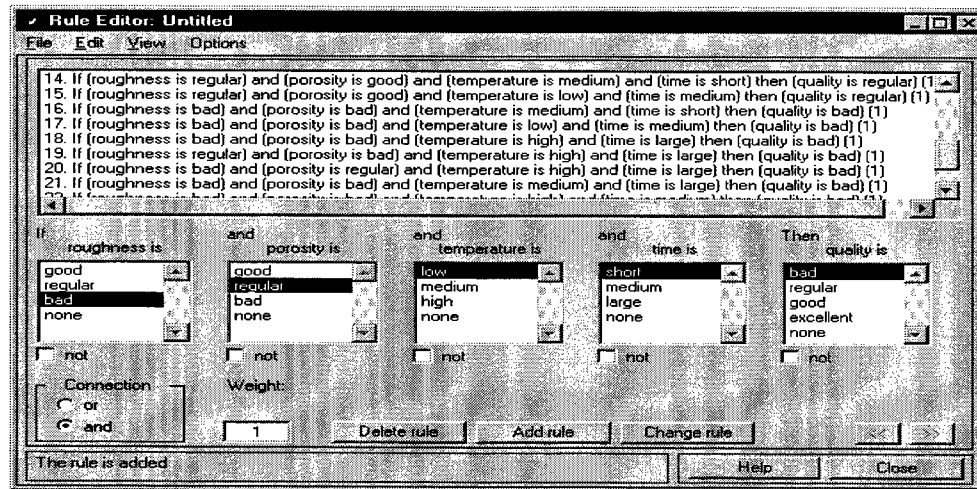


Fig. 3. Fuzzy rule base in the rule editor of the fuzzy logic toolbox of MATLAB.

Figure 4 shows the membership functions for the values of the "quality" linguistic variable. These functions were defined in the membership function editor of the fuzzy logic toolbox.

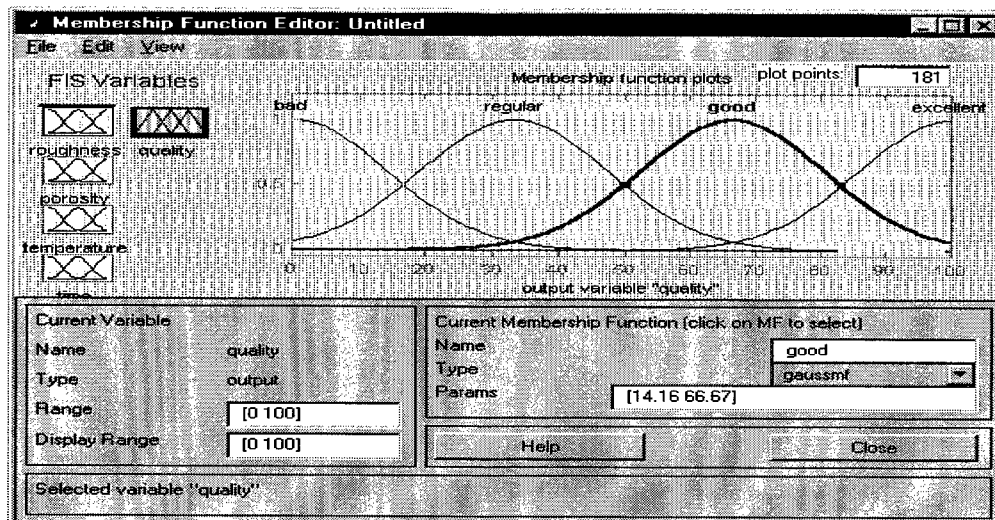


Fig. 4. Membership function plots for the linguistic values of the quality variable.

We show in Figure 5(a) the non-linear surface for the problem of quality evaluation using as input variables: roughness and porosity. We also show in Figure 5(b) the non-linear surface for the roughness and temperature. The three-dimensional surfaces represent the non-linear fuzzy model for the problem.

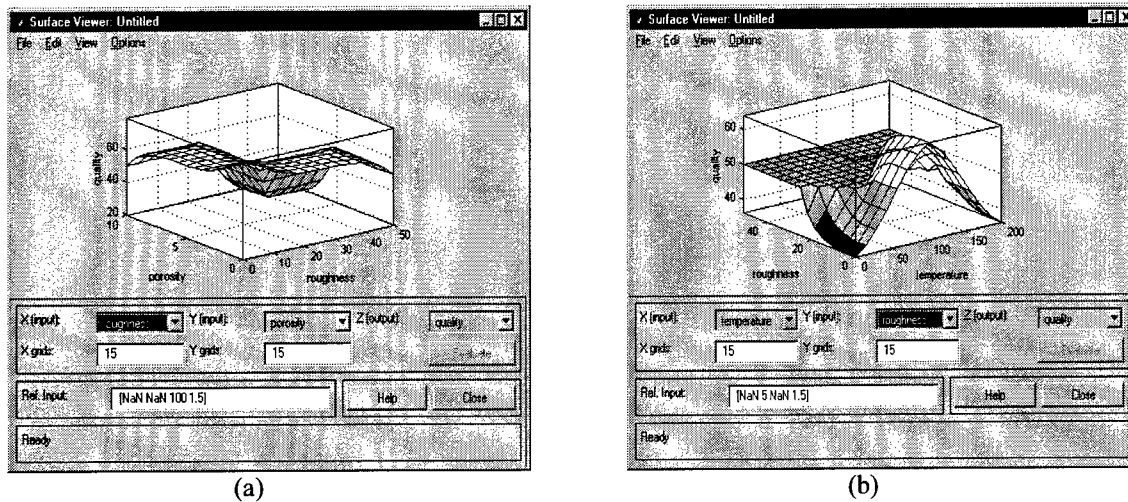


Fig. 5. Non-linear surface for quality evaluation with respect to
(a) roughness and porosity. (b) roughness and temperature.

We show in Figure 6 the reasoning procedure for quality control when specific values for the roughness and porosity are given. In this figure we can see how the final quality of the product is evaluated with the Mamdani inference system [4]. The results correspond to the values given by the real human experts for the domain of application.

We have to mention here that these simulation experiments for a specific problem of material processing show very good results. We have also tried our new fuzzy-fractal approach for quality control with other types of manufactured products with encouraging results.

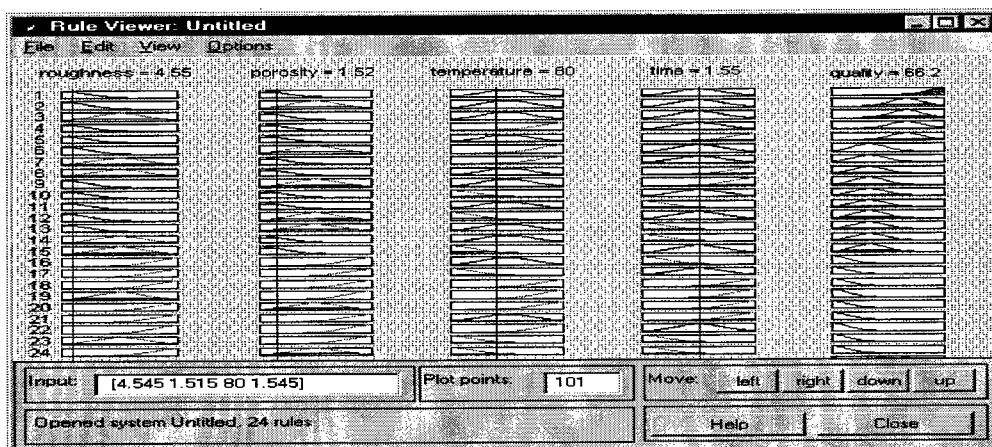


Fig. 6. Reasoning procedure for specific values of the roughness and porosity

CONCLUSIONS

We have very good simulation results in automated quality control, with our new fuzzy-fractal approach, for several types of manufactured products. For the specific case of material processing, the use of the fractal dimension to characterize the surface of the material is a good choice because it enables an efficient classification scheme for the roughness and porosity of the material. On the other hand, the use of fuzzy logic enables quality

evaluation of the product using as input the linguistic values of porosity and roughness. The new method for quality control combines the advantages of fuzzy logic (use of expert knowledge) with the advantage of fractal theory (an efficient classification scheme) to achieve the goal of automated quality control in material processing. Our new fuzzy-fractal approach can be used for different kinds of products because the fractal dimension can be used to characterize the geometrical complexity of the product, and fuzzy logic enables the simulation of the expert decision process for quality evaluation. We can conclude then that in this paper the authors have successfully generalized their previous work on fractal characterization of geometrical objects [1], with a computer system that can perform automated quality control of materials processing. Also, we can conclude that the combination of Soft Computing techniques and Fractal Theory is giving us a better method for quality control (more efficient in time and in accuracy).

REFERENCES

1. Castillo, O., Melin, P., 1994. Developing a New Method for the Identification of Microorganisms for the Food Industry using the Fractal Dimension, *Journal of Fractals*, , 2(3), 457-460.
2. Castillo, O., Melin, P., 1995. Automated Quality Control in the Food Industry Combining Artificial Intelligence Techniques and Fractal Theory, *Proceedings of the Tenth International Conference on Applications of artificial intelligence in Engineering*, Wessex Institute of Technology, U.K., 109-118.
3. Castillo, O., Melin, P., 1996. Automated Quality Control for Manufacturing in the Food Industry using Fuzzy Logic and Fractal Theory, *Proceedings of DKSME96 Conf.*, Univ. of Arizona, USA, 349-360.
4. Mamdani, E.H., Assilian, S., 1975. An Experiment in Linguistic Synthesis with a Fuzzy Logic Controller, *Man-Machine Studies*, 7, 1-13.
5. Mandelbrot, B., 1983. *The Fractal Geometry of Nature*, W. H. Freeman and Company.
6. Maurice, G., Pigeat, P., Chagroune, L., Weber, B., Thomas, A., 1994. Fractal Micro-roughness Model and Optical Emissivity Properties of Metallic Surfaces, *Fractals*, World Scientific, 2(2), 233-236.
7. Melin, P., Castillo O., 1996. Modelling and Simulation for Bacteria Growth Control in the Food Industry using Artificial Intelligence, *Proceedings CESA96*, Lille, France, 676-681.
8. Melin, P., Castillo O., 1997. Automated Mathematical Modelling and Simulation for Bacteria Growth Control in the Food Industry using Artificial Intelligence and Fractal Theory, *Systems, Analysis, Modelling and Simulation*, Gordon and Breach, 189-206.
9. Melin, P., Castillo, O., 1998. A New Method for Adaptive Model-Based Neuro-Fuzzy-Fractal Control of Non-Linear Dynamic Plants: The Case of Biochemical Reactors, *Proceedings of IPMU98*, EDK Publishers, Paris, France, 1, 475-482.
10. Nakamura, S., 1997. *Numerical Analysis and Graphic Visualization with MATLAB*, Prentice-Hall.
11. Yang, H.N., Wang, G.C., Lu, T.M., 1993. *Diffraction from Rough Surfaces and Dynamic Growth Fronts*, World Scientific, Singapore.
12. Zadeh, L.A., 1975. The Concept of a Linguistic Variable and its Application to Approximate Reasoning, *Information Sciences*, 8, 43-80.

A STUDY ON AXISYMMETRIC INDENTATION BY THE RIGID-PLASTIC FINITE-BOUNDARY ELEMENT METHOD

Y.-M. Guo* and K. Nakanishi*

** Department of Mechanical Engineering, Kagoshima University, 1-21-40 Korimoto,
Kagoshima City, 890-0065, Japan*

ABSTRACT

As compared with the conventional rigid-plastic finite element methods, the rigid-plastic finite-boundary element method is formulated with mixed-type. Therefore, this method possesses a merit in that this method can cover the compatibility of not only nodal velocity but also nodal velocity's derivative. On the other hand, the rigid-plastic finite-boundary element method does not need repetitional calculations in any computing step. Therefore, this method possesses another merit in that there is not any divergence possibility of repetitional calculations. An axisymmetric indentation problem is analyzed by the rigid-plastic finite-boundary element method in this paper. The processes from 0% to 30% reduction in the vertical height are simulated. Contours of effective strain, effective strain rate, effective stress and shear stress, etc. are obtained successfully.

INTRODUCTION

The rigid-plastic finite-boundary element method is formulated with mixed-type where the mixed variables are nodal velocities and derivatives of nodal velocity. Therefore, this method can cover the compatibility of not only nodal velocity but also nodal velocity's derivative, and nodal velocities and derivatives of nodal velocity can be calculated with the same precision for this method. While, the conventional rigid-plastic finite element methods are formulated with single-type, then these methods can not cover the compatibility of nodal velocity's derivative, and nodal velocities and derivatives of nodal velocity can not be calculated with the same precision. On the other hand, the rigid-plastic finite-boundary element method is a kind of solution in open form, then it does not need repetitional calculations in any computing step. Therefore, this method possesses another merit in that there is not any divergence possibility of repetitional calculations. While, the conventional rigid-plastic finite element methods are almost a kind of solution in closed form which needs repetitional calculations in every computing step, then there are some divergence possibilities of repetitional calculations in the conventional methods.

In this paper, the axisymmetric rigid-plastic finite-boundary element method is formulated, and an axisymmetric indentation problem is analyzed by this method.

AXISYMMETRIC FORMULATION

For a small axisymmetric element, the differential equations on mechanics equilibrium can be expressed as (In this paper the body forces are omitted for simplicity.):

$$\frac{\partial \sigma_R}{\partial R} + \frac{\partial \sigma_{RZ}}{\partial Z} + \frac{\sigma_R - \sigma_\theta}{R} = 0 \quad (1a)$$

$$\frac{\partial \sigma_{RZ}}{\partial R} + \frac{\partial \sigma_Z}{\partial Z} + \frac{\sigma_{RZ}}{R} = 0 \quad (1b)$$

It may be assumed that the small element consists of a kind of rigid-plastic material. By the theory of compressible plasticity[1] which introduces a term of dilatation into the yield criterion, the relation equation of stress and strain rate can be written as:

$$\begin{bmatrix} \sigma_R \\ \sigma_Z \\ \sigma_\theta \\ \sigma_{RZ} \end{bmatrix} = \frac{\sigma_{eq}}{\dot{\epsilon}_{eq}} \left\{ \frac{1}{3} \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} + \left(\frac{1}{g} - \frac{2}{9} \right) \begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \right\} \begin{bmatrix} \frac{\partial u}{\partial R} \\ \frac{\partial v}{\partial Z} \\ \frac{u}{R} \\ \frac{\partial u}{\partial Z} + \frac{\partial v}{\partial R} \end{bmatrix} \quad (2)$$

where g is a material constant that indicates dependence of hydrostatic stress. σ_{eq} and $\dot{\epsilon}_{eq}$ denote the effective stress and the effective strain rate, respectively. u and v denote velocity components.

Because we can assume the ratio of $\sigma_{eq}/\dot{\epsilon}_{eq}$ in the small element to be a constant. Substituting equation (2) into equations (1), the following equations of mechanics equilibrium on the small axisymmetric element are derived:

$$r^2 u + \left(\frac{1}{3} + \frac{3}{g} \right) \left(\frac{\partial^2 u}{\partial R^2} + \frac{\partial^2 v}{\partial R \partial Z} \right) + \frac{1}{R} \left(\frac{4}{3} + \frac{3}{g} \right) \left(\frac{\partial u}{\partial R} - \frac{u}{R} \right) = 0 \quad (3a)$$

$$r^2 v + \left(\frac{1}{3} + \frac{3}{g} \right) \left(\frac{\partial^2 u}{\partial R \partial Z} + \frac{\partial^2 v}{\partial Z^2} \right) + \frac{1}{R} \left[\left(\frac{1}{3} + \frac{3}{g} \right) \frac{\partial u}{\partial Z} + \frac{\partial v}{\partial R} \right] = 0 \quad (3b)$$

It may be considered that the rigid-plastic small element bears two load systems (forced system and observed system). Then, Samigiana's equation[2] in the case of that forced point lie outside analysis zone can be applied to the small element. Therefor, the following integral equations on the small axisymmetric element Ω are given:

$$\int_{\Omega} u^*(p, Q) r^2 u(Q) d\Omega(Q) + \int_{\Omega} u^*(p, Q) \left(\frac{1}{3} + \frac{3}{g} \right) \left[\frac{\partial^2 u(Q)}{\partial R^2(Q)} + \frac{\partial^2 v(Q)}{\partial R(Q) \partial Z(Q)} \right] d\Omega(Q) + \int_{\Omega} u^*(p, Q) \left(\frac{4}{3} + \frac{3}{g} \right) \frac{1}{R(Q)} \left[\frac{\partial u(Q)}{\partial R(Q)} - \frac{u(Q)}{R(Q)} \right] d\Omega(Q) = 0 \quad (4a)$$

$$\int_{\Omega} u^*(p, Q) r^2 v(Q) d\Omega(Q) + \int_{\Omega} u^*(p, Q) \left(\frac{1}{3} + \frac{3}{g} \right) \left[\frac{\partial^2 u(Q)}{\partial R(Q) \partial Z(Q)} + \frac{\partial^2 v(Q)}{\partial Z^2(Q)} \right] d\Omega(Q) + \int_{\Omega} u^*(p, Q) \frac{1}{R(Q)} \left[\left(\frac{1}{3} + \frac{3}{g} \right) \frac{\partial u(Q)}{\partial Z(Q)} + \frac{\partial v(Q)}{\partial R(Q)} \right] d\Omega(Q) = 0 \quad (4b)$$

where $u^*(p, Q)$ is a fundamental solution, which is a function of distance between a forced point p (lying outside analysis zone) and an observed point Q (lying inside analysis zone or on boundary of analysis zone):

$$u^*(p, Q) = \frac{1}{r(p, Q)} \quad (5)$$

where $r(p, Q)$ is the distance between point p and point Q . We can apply Green's equation to the first terms of equations (4), and obtain the following equations:

$$\int_{\Gamma} \left(\int_0^{2\pi} u^* d\theta \right) R q_r d\Gamma' - \int_{\Gamma} \left(\int_0^{2\pi} q^* d\theta \right) R u d\Gamma' + \int_{\Omega} \left(\int_0^{2\pi} u^* d\theta \right) \left(\frac{1}{3} + \frac{3}{g} \right) \left(\frac{\partial^2 u}{\partial R^2} + \frac{\partial^2 v}{\partial R \partial Z} \right) R dR dZ + \int_{\Omega} \left(\int_0^{2\pi} u^* d\theta \right) \left(\frac{4}{3} + \frac{3}{g} \right) \left(\frac{\partial u}{\partial R} - \frac{u}{R} \right) R dR dZ = 0 \quad (6a)$$

$$\int_{\Gamma} \left(\int_0^{2\pi} u^* d\theta \right) R q_z d\Gamma' - \int_{\Gamma} \left(\int_0^{2\pi} q^* d\theta \right) R v d\Gamma' + \int_{\Omega} \left(\int_0^{2\pi} u^* d\theta \right) \left(\frac{1}{3} + \frac{3}{g} \right) \left(\frac{\partial^2 u}{\partial R \partial Z} + \frac{\partial^2 v}{\partial Z^2} \right) R dR dZ + \int_{\Omega} \left(\int_0^{2\pi} u^* d\theta \right) \left[\left(\frac{1}{3} + \frac{3}{g} \right) \frac{\partial u}{\partial Z} + \frac{\partial v}{\partial R} \right] R dR dZ = 0 \quad (6b)$$

where

$$q_u = q_u(Q) = \frac{\partial u(Q)}{\partial n(Q)}; \quad q_v = q_v(Q) = \frac{\partial v(Q)}{\partial n(Q)} \quad (7)$$

$$\dot{q} = \dot{q}(p, Q) = \frac{\partial \dot{u}(p, Q)}{\partial n(Q)} = -\frac{1}{r^2(p, Q)} \frac{\partial r(p, Q)}{\partial n(Q)} \quad (8)$$

where Γ' and Ω' are the boundary and the domain on meridional plane of the small axisymmetric element Ω respectively. We do a discretization for the small axisymmetric element, that is dissociating the boundary Γ' into L boundary elements and the domain Ω' into one finite element. Then, the following linear equations that variables are nodal velocities and nodal velocity's normal change rates are obtained:

$$\sum_{i=1}^L \int_{\Gamma'_i} \bar{u}^* R N_{i,u} q_{u,i} d\Gamma' - \sum_{i=1}^L \int_{\Gamma'_i} \bar{q}^* R N_{i,u} u_i d\Gamma' + \int_{\Omega'} \bar{u}^* \left(\frac{1}{3} + \frac{3}{g} \right) \left(\frac{\partial^2 N u_u}{\partial R^2} + \frac{\partial^2 N v_u}{\partial R \partial Z} \right) R dR dZ + \int_{\Omega'} \bar{u}^* \left(\frac{4}{3} + \frac{3}{g} \right) \left(\frac{\partial N u_u}{\partial R} - \frac{N u_u}{R} \right) R dR dZ = 0 \quad (9a)$$

$$\sum_{i=1}^L \int_{\Gamma'_i} \bar{u}^* R N_{i,u} q_{u,i} d\Gamma' - \sum_{i=1}^L \int_{\Gamma'_i} \bar{q}^* R N_{i,v} v_i d\Gamma' + \int_{\Omega'} \bar{u}^* \left(\frac{1}{3} + \frac{3}{g} \right) \left(\frac{\partial^2 N u_u}{\partial R \partial Z} + \frac{\partial^2 N v_u}{\partial Z^2} \right) R dR dZ + \int_{\Omega'} \bar{u}^* \left[\left(\frac{1}{3} + \frac{3}{g} \right) \frac{\partial N u_u}{\partial Z} + \frac{\partial N v_u}{\partial R} \right] R dR dZ = 0 \quad (9b)$$

where u_u and v_u are nodal velocities, $q_{u,i}$ and $q_{v,i}$ are nodal velocity's normal change rates. $N_{i,u}$ and $N_{i,v}$ are the shape function vector of boundary element and that of finite element, respectively.

$$\bar{u}^* = \bar{u}^*(p, Q) = \int_0^{2\pi} u^*(p, Q) d\theta(Q) \quad (10)$$

$$\bar{q}^* = \bar{q}^*(p, Q) = \int_0^{2\pi} q^*(p, Q) d\theta(Q) \quad (11)$$

In equation (9), the forced points p lie outside the finite element and the observed points Q lie at the boundary element. An axisymmetric workpiece can be dissociated into M small axisymmetric elements, then linear equations on the axisymmetric workpiece can be obtained by globalization.

$$\sum_{m=1}^M \sum_{i=1}^L \int_{\Gamma'_i} \bar{u}^* R N_{i,u} q_{u,i} d\Gamma' - \sum_{m=1}^M \sum_{i=1}^L \int_{\Gamma'_i} \bar{q}^* R N_{i,u} u_i d\Gamma' + \sum_{m=1}^M \int_{\Omega'_m} \bar{u}^* \left(\frac{1}{3} + \frac{3}{g} \right) \left(\frac{\partial^2 N u_u}{\partial R^2} + \frac{\partial^2 N v_u}{\partial R \partial Z} \right) R dR dZ + \sum_{m=1}^M \int_{\Omega'_m} \bar{u}^* \left(\frac{4}{3} + \frac{3}{g} \right) \left(\frac{\partial N u_u}{\partial R} - \frac{N u_u}{R} \right) R dR dZ = 0 \quad (12a)$$

$$\sum_{m=1}^M \sum_{i=1}^L \int_{\Gamma'_i} \bar{u}^* R N_{i,u} q_{u,i} d\Gamma' - \sum_{m=1}^M \sum_{i=1}^L \int_{\Gamma'_i} \bar{q}^* R N_{i,v} v_i d\Gamma' + \sum_{m=1}^M \int_{\Omega'_m} \bar{u}^* \left(\frac{1}{3} + \frac{3}{g} \right) \left(\frac{\partial^2 N u_u}{\partial R \partial Z} + \frac{\partial^2 N v_u}{\partial Z^2} \right) R dR dZ + \sum_{m=1}^M \int_{\Omega'_m} \bar{u}^* \left[\left(\frac{1}{3} + \frac{3}{g} \right) \frac{\partial N u_u}{\partial Z} + \frac{\partial N v_u}{\partial R} \right] R dR dZ = 0 \quad (12b)$$

Number of the linear equations on the workpiece is determined on number of forced points p . The distance between every two forced points p ought to be longer than about 0.2 mm to keep independence of every equation.

\bar{u}^* and \bar{q}^* can be calculated as the following equations:

$$\bar{u}^* = \frac{4K(s)}{\sqrt{r_1^2(p, Q) + 4|R(p)|R(Q)}} \quad \text{if } \gamma > \gamma_0 \quad (13a)$$

$$\bar{u}^* = \sqrt{\frac{4}{|R(p)|R(Q)}} Q_{1/2}(\gamma) \quad \text{if } \gamma \leq \gamma_0 \quad (13b)$$

$$\bar{q} = \frac{4}{\sqrt{r_1^2(p,Q) + 4|R(p)|R(Q)}} \left\{ \frac{1}{2R(Q)} \left[\frac{R^2(p) - R^2(Q) + [Z(p) - Z(Q)]^2}{r_1^2(p,Q)} E(s) - K(s) \right] \frac{\partial R(Q)}{\partial n(Q)} + \frac{Z(p) - Z(Q)}{r_1^2(p,Q)} E(s) \frac{\partial Z(Q)}{\partial n(Q)} \right\} \quad \text{if } \gamma > \gamma_0 \quad (14a)$$

$$\bar{q} = - \sqrt{\frac{4}{|R(p)|R(Q)}} \frac{1}{R(Q)} \left\{ \frac{1}{2} \left[Q_{-1/2}(\gamma) + \frac{R^2(p) - R^2(Q) + [Z(p) - Z(Q)]^2}{|R(p)|R(Q)} \frac{dQ_{-1/2}(\gamma)}{d\gamma} \right] \frac{\partial R(Q)}{\partial n(Q)} + \frac{Z(p) - Z(Q)}{|R(p)|} \frac{dQ_{-1/2}(\gamma)}{d\gamma} \frac{\partial Z(Q)}{\partial n(Q)} \right\} \quad \text{if } \gamma \leq \gamma_0 \quad (14b)$$

where

$$r_1(p,Q) = \sqrt{[R(p) - R(Q)]^2 + [Z(p) - Z(Q)]^2} \quad (15)$$

$$s = s(p,Q) = \frac{4|R(p)|R(Q)}{r_1^2(p,Q) + 4|R(p)|R(Q)} \quad 0 < s < 1 \quad (16)$$

$$\gamma = \gamma(p,Q) = 1 + \frac{r_1^2(p,Q)}{2|R(p)|R(Q)} \quad \gamma > 1 \quad (17)$$

$\gamma_0 \approx 1.05 \sim 1.10$

$K(s)$ denotes the 1st kind of full elliptic integration, and can be approximately calculated as the following equation if $s < 1$:

$$K(s) = \frac{\pi}{2} \left(1 + \frac{1^2}{2^2} s^2 + \frac{1^2 \cdot 3^2}{2^2 \cdot 4^2} s^4 + \frac{1^2 \cdot 3^2 \cdot 5^2}{2^2 \cdot 4^2 \cdot 6^2} s^6 + \dots \right) \quad (18)$$

$E(s)$ denotes the 2nd kind of full elliptic integration, and can be approximately calculated as the following equation if $s < 1$:

$$E(s) = \frac{\pi}{2} \left(1 - \frac{1}{2^2} s^2 - \frac{1^2 \cdot 3}{2^2 \cdot 4^2} s^4 - \frac{1^2 \cdot 3^2 \cdot 5}{2^2 \cdot 4^2 \cdot 6^2} s^6 - \dots \right) \quad (19)$$

$Q_{-1/2}(\gamma)$ denotes the 2nd kind of Legendre function, and can be calculated as the following equations[2] if γ is small:

$$Q_{-1/2}(\gamma) = -\frac{1}{2} \ln \left(\frac{\gamma-1}{32} \right) \quad (20)$$

$$\frac{dQ_{-1/2}(\gamma)}{d\gamma} = -\frac{1}{2(\gamma-1)} \quad (21)$$

To avoid singularity, the following limitations are adopted:

$$R(Q) \geq 0.1 \text{ [mm]} \quad (22)$$

$$|R(p)| \geq 0.1 \text{ [mm]} \quad (23)$$

$$r_1(p,Q) \geq 0.2 \text{ [mm]} \quad (24)$$

Nodal boundary conditions can be written as:

1) Velocity boundary condition:

$$u_\alpha = u_c; \quad v_\alpha = v_c \quad (25)$$

2) Friction boundary condition:

$$p_t = \mu p_n \quad (26)$$

3) Free surface boundary condition:

$$p_n = 0; \quad p_t = 0 \quad (27)$$

where u_c and v_c are known nodal velocity vectors. μ is a friction factor. p_n and p_t are a nodal normal pressure and a nodal tangent pressure, respectively, which can be calculated from nodal stress vector σ_α . For every node, number of the nodal variables must be equal to number of nodal equations (including the number of nodal boundary conditions).

RESULTS OF ANALYSIS

The object of analysis is an axisymmetric indentation, only the upper half of the workpiece being considered owing to symmetry. The processes from 0% to 30% reduction in vertical height of dies are simulated. The three-noded quadratic boundary element and the eight-noded quadratic finite element are adopted on the meridional plane of the axisymmetric workpiece. L in equations (9) and (12) is taken as 4. Whilst g is assumed as 0.01. Friction factor μ is taken as 0.1. The velocity of dies is 1.0 mm/s.

Figure 1 depicts contours of effective strain rate at 10% reduction (where one quadrant of the workpiece is shown). It may be seen that the effective strain rate near to the centre zone is large, whilst those near to the dies and the outer middle zone are small.

Figure 2 and figure 3 depict contours of effective strain, in which figure 2 is the calculated result at 20% reduction and figure 3 is that at 30% reduction. As seen in these figures, the effective strain near to the centre zone is large, whilst that of the outer zone is small at the two reductions, and values of the effective strain follow increase in the reduction.

Figure 4 depicts contours of effective stress at 20% reduction. Comparing this figure and figure 2, it is seen that the distributions of effective stress and effective strain are about the same. It is because that the effective stress is only a function of the effective strain for this paper.

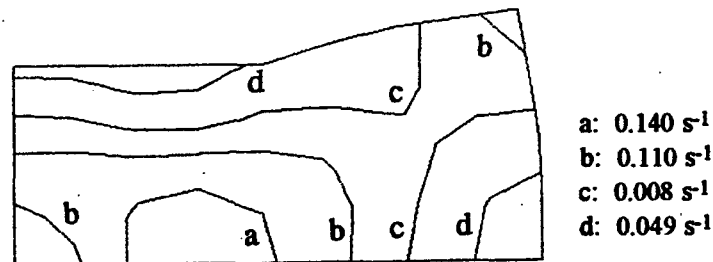


Fig. 1. Contours of effective strain rate at 10% reduction.

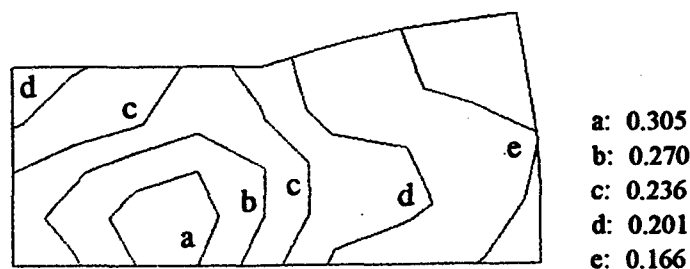


Fig. 2. Contours of effective strain at 20% reduction.

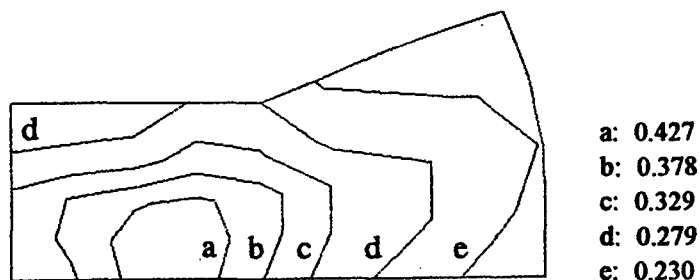


Fig. 3. Contours of effective strain at 30% reduction.

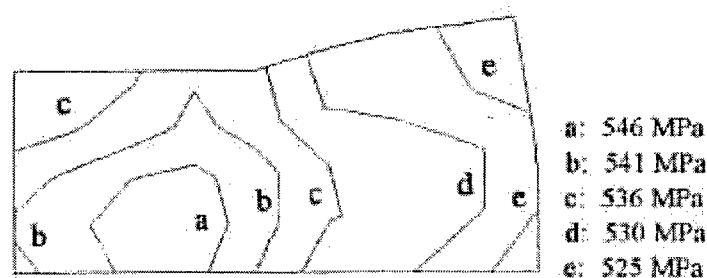


Fig. 4. Contours of effective stress at 20% reduction.

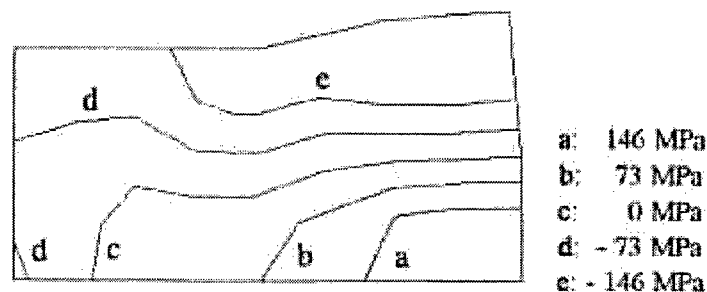


Fig. 5. Contours of shear stress at 10% reduction.

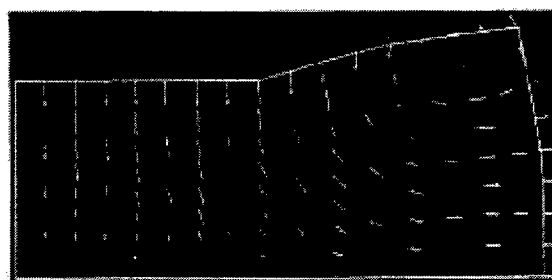


Fig. 6. velocity fields at 20% reduction.

Figure 5 shows contours of shear stress at 10% reduction. As seen in this figure, the shear stresses of the outer end zones and outer middle zone are large, whilst that among the zones of large shear stress is small.

Figure 6 is velocity field at 20% reduction. The flow pattern of the material can be seen from this figure.

CONCLUSIONS

The rigid-plastic finite-boundary element method can cover the compatibility of not only nodal velocity but also nodal velocity's derivative, and does not need repetitional calculations in any computing step. Therefore it is considered that this method is a better method than the conventional rigid-plastic finite element methods. The equations in axisymmetric form of the rigid-plastic finite-boundary element method have been formulated, and the axisymmetric indentation process has been simulated.

REFERENCES

1. K. Mori, K. Osakada and T. Oda, 1982. Simulation of plane strain rolling by the rigid-plastic finite element method. *Int. J. Mech. Sci.*, 24, 519-527.
2. C. A. Brebbia, J. C. F. Telles and L. C. Wrobel, 1984. Boundary element techniques—theory and applications in engineering, Springer-Verlag.

Design of Intelligent Spindle for High Speed Machining

B.L.Zhang, Y.P.Li, B.S.Zhu, P.Ma and Y.Luo

Dept. of Mechanical and Electronic Engineering
Guangdong University of Technology
729 East Dongfeng Road, Guangzhou 510090, CHINA
Tel/Fax: 86-20-87613563 Email: blzhang@gdut.edu.cn

ABSTRACT

The Spindle is the most important component in a high speed CNC machine tool. Its thermal and dynamic behaviors are almost decided by the drive system and the support system of the spindle. Some special requirements for high speed spindles must be satisfied. C-axis control is also required. A high speed, high power spindle driven by a built-in motor has been developed. The frameless spindle motor is located between the bearings. The frameless motor is cooled by a air-water cooling system. Si_3N_4 ceramic ball bearings are used to support the spindle. An oil-air lubricator is used to reduce the friction force and control the temperature-rise. Some design problems are discussed for improving the behaviour of the machine tool.

Robotics and Intelligent Control I

Autonomous Control of Complex Dynamical Systems in Support of a Manned Mission to Mars

James A. Kurien, Daniel J. Clancy

NASA Ames Research Center
MS 269-2, Moffett Field, CA 94035

Email: kurien@ptolemy.arc.nasa.gov, clancy@ptolemy.arc.nasa.gov

ABSTRACT

Space missions have historically relied upon a large ground staff, numbering in the hundreds for complex missions, to maintain routine operations. When an anomaly occurs, this small army of engineers attempts to identify and work around the problem. A piloted Mars mission, with its multiyear duration, cost pressures, half-hour communication delays and two-week blackouts cannot be closely controlled by a battalion of engineers on Earth. Flight crew involvement in routine system operations must also be minimized to maximize science return. It also may be unrealistic to require the crew have the expertise in each mission subsystem needed to diagnose a system failure and effect a timely repair, as engineers did for Apollo 13.

Enter model-based autonomy, which allows complex systems to autonomously maintain operation despite failures or anomalous conditions, contributing to safe, robust, and minimally supervised operation of spacecraft, life support, ISRU and power systems. Autonomous reasoning is central to the approach. A reasoning algorithm uses a logical or mathematical model of a system to infer how to operate the system, diagnose failures and generate appropriate behavior to repair or reconfigure the system in response.

The "plug-and-play" nature of the models enables low cost development of autonomy for multiple platforms. Declarative, reusable models capture relevant aspects of the behavior of simple devices (e.g. valves or thrusters). Reasoning algorithms combine device models to create a model of the system-wide interactions and behavior of a complex, unique artifact such as a spacecraft. Rather than requiring engineers to envision all possible interactions and failures at design time or perform analysis during the mission, the reasoning engine generates the appropriate response to the current situation, taking into account its system-wide knowledge, the current state, and even sensor failures or unexpected behavior.

INTRODUCTION

Exploring and ultimately settling Mars will be a milestone in the development of our civilization and an uncompromising measure of our courage, cleverness and resolve. Accordingly, it will also be an unprecedented technical challenge, involving multiple interdependent mission elements, multiyear duration, incredible budgetary pressure and the duty to protect human lives in a harsh environment millions of miles from Earth. Evidence of the utility of highly capable, robust and coordinated autonomous systems in meeting this challenge pervades mission scenarios such as Mars Direct [1] and the NASA Mars Reference Mission [2].

Model-based autonomy involves the use of automated reasoning engines and high level models of the system being controlled to generate correct system behavior on the fly, even in the face of failures or anomalous situations. This approach is proving to be a robust and cost effective method for developing more highly capable autonomous systems than have been deployed in the past and might prove invaluable to the development of piloted missions to Mars.

The next section of this paper describes why autonomous systems are needed to explore Mars. Section 3 discusses how this work can contribute to cheap, safe, robust, and minimally supervised systems on Mars. Section 4 describes Livingstone, one of the reasoning engines developed at Ames that will be tested onboard a spacecraft next year. Section 5 describes a number of Mars-related testbeds that are making use of model-based autonomy technology.

THE UTILITY OF AUTONOMY ON MARS

The need for robust, inexpensive and productive operation of remote assets on Mars appears throughout both the Mars Direct scenario and the Mars Reference Mission. In both of these mission designs, initial mission elements such as in-situ propellant production (ISPP) plants and the crew return vehicle must be able to operate for a period years in a harsh environment with limited downlink capabilities and a reduced set of ground control personnel. Such systems must maintain efficient operation in spite of unexpected failures, novel environmental phenomena and degraded system capabilities. Safety places high demands on system robustness: the crew cannot depart Earth if propellant plant down time results in inadequate production or if the return vehicle cannot verify nominal operation.

Once the crew does depart Earth, they will be travelling two orders of magnitude farther from home than the Apollo crews. They will be separated from mission control by thirty-minute communication delays and potentially multi-day communication blackouts imposed by the relative positions of Mars and the Earth. There will be a number of systems upon which the crew's ability to reach Mars or survive an abort to Earth will depend: life support, attitude control, propulsion, communications and power generation are examples. While only life support might seem to require immediate response to anomalies, many other situations require on board response as well: losing attitude control during an aerobreaking maneuver, failures which need to be quickly saved, and loss of communications with Earth are all cases in point.

Once on the Martian surface, maximization of exploration becomes a focus in addition to safety. We do not yet have the resources to send crews of fifteen to Mars to run a Martian science outpost and support system. Hence crew involvement in routine operations such as controlling the life support system or maintaining rovers must be minimized and minor anomalies must be resolved locally rather than awaiting ground analysis. In addition, to maximize science return in this unknown environment, operations on Mars must be able to rapidly adapt to take advantage of new science opportunities or make the best of degraded capabilities.

These challenges to maintaining safety and productivity on Mars from Earth for several years are daunting when one considers the current state of mission operations. Current piloted missions rely upon near-instantaneous contact with hundreds of engineers and operators on the ground. In addition, recent attempts to teleoperate relatively simple systems for ninety days on Mars resulted in a considerable fraction of the mission being used to determine the state of the remote system and return it to productive operation, often over the course of a day or more [3].

The Reference Mission therefore explicitly calls for autonomous systems on Mars to allow unmanned systems to robustly prepare for human arrival, to protect crew and resources by rapidly responding to critical failures, to free explorers from routine operations and to control operations costs for this complex, multi-year mission. In this context, autonomy means the ability to correctly react to a wide range of circumstances, both usual and anomalous, without the need for direct human supervision. If available, a robust onboard autonomy capability would enable safer, more affordable missions to Mars by allowing complex systems such as life support systems or spacecraft attitude control systems to operate for extended periods of time without supervision over a wide range of nominal and anomalous operating conditions. The benefits would be increased safety and reduced downtime for mission critical systems, leverage of scarce human skills by automation of routine tasks, and reduced ground operations due to unattended recovery from anomalies and less detailed commanding requirements.

Currently, NASA's operational experience with the type of high capability, failure-tolerant autonomy described in the Reference Mission is low. To date, no fully automated power plants, life support, or cryogen plants have been deployed. Some automated planning and scheduling has been used to pre-compute command sequences for spacecraft and to schedule space shuttle refurbishment, but no deployed system has autonomously replanned its mission activities in the field. In addition, the robotic systems that have been deployed in space have been almost entirely dependent upon pre-computed command sequences relayed from Earth controllers, and have not been highly autonomous in the sense conveyed above.

Of course, every unmanned system sent into space has required some level of autonomy: if a spacecraft cannot at least point its antenna at Earth and wait for help after the expected kinds of failures, it is likely to

be lost. Currently programmers and mission control operators use their commonsense understanding of hardware and mission goals to produce code and control sequences that will allow a spacecraft or other system to achieve some goal while allowing for some (usually very small) amount of uncertainty in the environment. This has the disadvantages of being relatively time intensive, error prone, and not particularly reusable. Because of the amount of analysis involved, such systems usually allow for uncertainty by being extremely conservative and provide the minimal amount of adaptability necessary to raise the likelihood of survival of the spacecraft. If an anomaly occurs the spacecraft or other system typically halts all activity, achieves a safe mode, and awaits further instructions. One notable exception is the attitude and articulation control system on the Cassini spacecraft, which represents state-of-the-art in deployed spacecraft autonomy [4] and which has not been replicated on the "faster, better, cheaper" missions which have followed.

The cost to develop highly robust autonomous control software and the ability of such systems to improve safety and productivity of assets deployed on Mars (or deep space or Europa for that matter) are significant risk factors that impact NASA's ability to accurately plan and scope future missions. One intent of the work described in the paper, model-based autonomy, is to demonstrate that highly robust autonomous systems can be developed more easily and more cheaply than the more modest systems which have been deployed to date.

What is model-based autonomy?

Model-based autonomy refers to the achievement of robust, autonomous operation through a growing set of reusable artificial intelligence (AI) reasoning algorithms that reason about first principles models of physical systems (e.g. spacecraft). In this context, a *model* is a logical or mathematical representation of a physical object or piece of software. A *first principles* model captures what is true about behavior or structure of the object (e.g. fluid flows through an open valve unless it is clogged). This is as opposed to traditional programs or rule-based expert systems, which capture what to do (e.g. turn on valves A, B, & C to start fuel flow) but unfortunately work only in certain implicit contexts (e.g. valves A, B, & C are working and A, B & C happen to control the fuel flow).

Since model-based autonomous systems do not contain an encoding of what to do in each situation, they must reason about the appropriate action to take or conclusion to draw based upon their models and the currently available information about the environment. The past few decades of AI research have produced reasoning engines that can plan a course of action to achieve a goal, identify the current state of a physical system, reconfigure that system to enable some function (e.g. make the engine thrust) and so on from a first principles model.

Reasoning directly from the model, the current observations of the world and the task at hand provides many large advantages over traditional software development. Not least among these are that the system is robust in uncertain environments since it was not hard coded to respond to certain situations, the models and inference engines can be reused, and the models explicitly capture the assumptions about the system that are being relied upon to control it.

Benefits of Model-Based Autonomy

A model-based autonomous controller provides a number of benefits that are critical in the development of a robust control architecture required to support a manned mission to Mars. While various control techniques have been developed over the years, many of these techniques focus on the low-level control response required to maintain the system within a stable operating regime. The control techniques are then augmented with a higher-level discrete controller that is often implemented using a traditional software development methodology. Developing controllers in this manner is often time-consuming and often the resulting controllers are limited in their ability to handle novel component interactions that were not explicitly anticipated by the software developer. Furthermore, for devices that must operate without operator intervention for extended periods of time, it is often quite difficult if not impossible to write software that can handle all of the possible combinations of faults that can potentially occur over time.

A model-based controller addresses this problem by using a declarative specification of the device being controlled and the goals that are to be achieved. This results in increased safety and reliability while providing a significant decrease in the overall development costs due to code reuse and the compositional

nature of the modeling paradigm. Finally, the use of a declarative specification simplifies the development of an advanced user interface that can be used to monitor the state of the device and to query the controller to obtain information about the actions that have been taken and the justification for these actions. This capability tends to decrease the overall system operation costs since the man power required to monitor the device is significantly reduced.

LIVINGSTONE

As mentioned previously, Livingstone is a model-based discrete controller. Its function is to infer the current state (mode) of each relevant device making up the system being controlled and to recommend actions that can reconfigure the system so that it achieves the currently desired configuration goals, if possible. In practice, these configuration goals could be provided by a human operating some apparatus by issuing high level configuration commands, or by some automated system such as the Smart Executive (Exec) mentioned above, which decomposes a high level plan into a series of configuration goals to be achieved. Purely for the sake of the discussion below, we will assume the Exec is providing the configuration goals and that the system being controlled is a spacecraft.

To track the modes of system devices, Livingstone eavesdrops on commands that are sent to the spacecraft hardware by the Exec. As each command is executed, Livingstone receives observations from spacecraft's sensors, abstracted by monitors in the real time control software for the Attitude Control Subsystem (ACS), communications bus, or whatever hardware is present. Livingstone combines these commands and observations with declarative models of the spacecraft components to determine the current state of the system and report it to the Exec. A pathologically simple example is shown schematically in Figure 1. In the nominal case, Livingstone merely confirms that the commands had the expected effect on spacecraft state. In case of failure, Livingstone diagnoses the failure and the current state of the spacecraft and provides a recovery recommendation. A single set of models and algorithms are exploited for command confirmation, diagnosis and recovery.

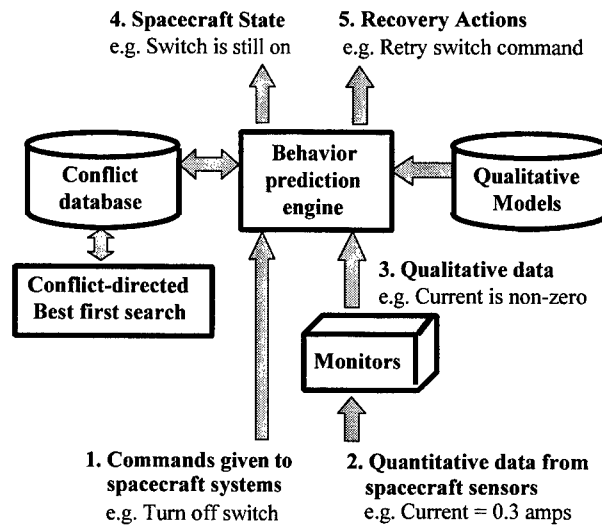


Fig. 1. Information Flow in Livingstone.

The capabilities of the Livingstone inference engine can be divided into two parts: mode identification (MI) and mode reconfiguration (MR). MI is responsible for identifying the current operating or failure mode of each component in the spacecraft. Following a component failure, MR is responsible for suggesting reconfiguration actions that restore the spacecraft to a configuration that achieves all current configuration goals required by the Exec. Livingstone can be viewed as a discrete model-based controller in which MI provides the sensing component and MR provides the actuation component. MI's mode inference allows the Exec to reason about the state of the spacecraft in terms of component modes or even high level capabilities such as "able to produce thrust" rather than in terms of low level sensor values. MR supports

the run-time generation of novel reconfiguration actions to return components to the desired mode or to re-enable high level capabilities such as "able to produce thrust".

Livingstone uses algorithms adapted from model-based diagnosis [11,12] to provide the above functions. The key idea underlying model-based diagnosis is that a combination of component modes is a possible description of the current state of the spacecraft only if the set of models associated with these modes is consistent with the observed sensor values. Following de Kleer and Williams [13], MI uses a conflict directed best-first search to find the most likely combination of component modes consistent with the observations. Analogously, MR uses the same search to find the least-cost combination of commands that achieve the desired goals in the next state. Furthermore, both MI and MR use the same system model to perform their function. The combination of a single search algorithm with a single model, and the process of exercising these through multiple uses, contributes significantly to the robustness of the complete system. Note that this methodology is independent of the actual set of available sensors and commands. Furthermore, it does not require that all aspects of the spacecraft state are directly observable, providing an elegant solution to the problem of limited observability.

The use of model-based diagnosis algorithms immediately provides Livingstone with a number of additional features. First, the search algorithms are sound and complete, providing a guarantee of coverage with respect to the models used. Second, the model building methodology is modular, which simplifies model construction and maintenance, and supports reuse. Third, the algorithms extend smoothly to handling multiple faults and recoveries that involve multiple commands. Fourth, while the algorithms do not require explicit fault models for each component, they can easily exploit available fault models to find likely failures and possible recoveries.

Livingstone extends the basic ideas of model-based diagnosis by modeling each component as a finite state machine, and the whole spacecraft as a set of concurrent, synchronous state machines. Modeling the spacecraft as a concurrent machine allows Livingstone to effectively track concurrent state changes caused either by deliberate command or by component failures. An important feature is that the behavior of each component state or mode is captured using abstract, or qualitative, models [14]. These models describe qualities of the spacecraft's structure or behavior without the detail needed for precise numerical prediction, making abstract models much easier to acquire and verify than quantitative engineering models. Examples of qualities captured are the power, data and hydraulic connectivity of spacecraft components and the directions in which each thruster provides torque. While such models cannot quantify how the spacecraft would perform with a failed thruster for example, they can be used to infer which thrusters are failed given only the signs of the errors in spacecraft orientation. Such inferences are robust since small changes in the underlying parameters do not affect the abstract behavior of the spacecraft. In addition, abstract models can be reduced to a set of clauses in propositional logic. This form allows behavior prediction to take place via unit propagation, a restricted and very efficient inference procedure.

It is important to note that the Livingstone models are not required to be explicit or complete with respect to the actual physical components. Often models do not explicitly represent the cause for a given behavior in terms of a component's physical structure. For example, there are numerous causes for a stuck switch: the driver has failed, excessive current has welded it shut, and so on. If the observable behavior and recovery for all causes of a stuck switch are the same, Livingstone need not closely model the physical structure responsible for these fine distinctions. Models are always incomplete in that they have an explicit unknown failure mode. Any component behavior that is inconsistent with all known nominal and failure modes is consistent with the unknown failure mode. In this way, Livingstone can still infer that a component has failed, though the failure was not foreseen or was simply left unmodeled because no recovery is possible. By modeling only to the level of detail required to make relevant distinctions in diagnosis (distinctions that prescribe different recoveries or different operation of the system) we can describe a system with qualitative "common-sense" models which are compact and quite easily written.

MARS RELATED APPLICATIONS

The intent behind model-based autonomy is to create generic, high capability reasoning systems that can be adapted to a wide range of applications simply by writing the appropriate models. As such, model-based

autonomy might be able to contribute to the control of a variety of elements of a piloted Mars mission. In this early stage of Mars mission definition, model-based autonomy is involved in the prototyping of a number of specific mission elements.

Closed-Loop Ecological Life Support Systems (CELSS)

In order to transport and support humans for Mars expeditions, NASA's Human Exploration and Development of Space (HEDS) requirements state a need for autonomous operation of life support, ISRU and transport equipment. During a Mars expedition, autonomous plant operations would allow unmanned systems to prepare for human arrival, protect crew and resources by rapidly responding to critical failures, and free humans from routine operations, allowing greater exploration.

At NASA's Johnson Space Center (JSC), a closed loop life support testbed called Bioplex has been constructed. The Bioplex consists of three sections: a three story cylindrical living quarters similar to the Mars habitats discussed in various mission proposals; a plant chamber where wheat is grown to provide food and exchange CO_2 for O_2 ; and an incinerator chamber used to eliminate solid waste and produce CO_2 . The most recent Bioplex testing is referred to as the Product Gas Transfer phase as it concentrates on generation and distribution of product gases (CO_2 from the crew and incinerator and O_2 from the plants) and does not yet address issues such as waste water recycling or power management.

A JSC advanced development group has developed an autonomous control system to operate the product gas transfer phase of Bioplex [15]. This system, based upon the 3T autonomy architecture [16], maintains the appropriate atmosphere in each chamber by extracting and storing product gases and coordinating activities such as firing the incinerator or opening the plant chamber for human access. The system successfully controlled gas transfer during test in which a human crew inhabited the Bioplex for ninety days. It was not expected, however, to maintain operation in the face of failures, though many would likely occur over a 4-year mission.

We are currently working to integrate the Livingstone mode identification and reconfiguration engine with JSC's 3T architecture, adding to it the ability to determine the current state of the testbed and respond to anomalous situations or failures by performing high level, system-wide reasoning. This will result in a single, reusable architecture which maintains the best possible operation of a regenerative life support system and other complex physical plants during both nominal operation and failures, somewhat analogous to the autonomic and immune functions of a living organism.

We intend to demonstrate the combined system by maintaining operation of the testbed over an extended test period and providing both fully autonomous and human-centered operation. To test the system, an outside examiner will be employed to introduce failures into the testbed as desired which the system will diagnose and attempt to mitigate.

The second goal is to demonstrate and extend the ability of model-based systems to reduce analysis, development and operations costs. The testbed application will be rapidly developed with tools that could be used to develop mission applications. Users will develop and operate the testbed by manipulating explicit models with visual tools. If previous experience is to be believed, far less effort will be required to develop, understand and revise the system than in an approach where system model is implicit but still must be maintained.

If successful, this demonstration will increase the likelihood that autonomy technologies being developed by NASA are appropriate and sufficiently mature when they are required for HEDS missions to Mars and other destinations. It will also ensure that the necessary technologies can be integrated and will identify needed extensions before such shortcomings could impact the critical path of a mission. In addition, JSC will have a prototype of a reusable, fault-tolerant, high-capability autonomous control system and the expertise to apply this system to a flight experiment or mission. This could be applied to any complex physical system that must be controlled and maintained over an extended period of time such as spacecraft, power plants, ISRU machinery, and autonomous or semiautonomous surface vehicles.

In-situ Resource Utilization

In-situ resource utilization, or "living off the land", is critical to making a piloted Mars mission robust and affordable [1]. More specifically, it is envisioned that in-situ propellant production (ISPP) plants will arrive on Mars years before humans and begin combining hydrogen brought from Earth with CO₂ from the Martian atmosphere to create methane. This fuel will power the ascent vehicle that will lift the crew off Mars to begin their trip home in addition to powering any methane-fueled surface vehicles the astronauts might possess.

Though the chemical reactions involved are conceptually quite simple, on Mars they are somewhat complicated by issues such as the low atmospheric pressure and slow contamination of the ISPP catalysts by trace elements in the Martian atmosphere. To ensure that adequate ISPP capability is available for future Mars missions, NASA has begun to explore ISPP designs and build prototype hardware for operation in Mars-like test chambers. Both JSC and NASA Kennedy Space Center (KSC) are involved in early ISPP development, and the KSC team is integrating Livingstone into their ISPP prototyping efforts.

The short-term focus of this collaboration is to integrate Livingstone's ability to diagnose and mitigate failures with existing KSC model-based technology to gain experience with a model-based monitoring, diagnosis and recovery system for ISPP. A secondary short-term goal is to determine if any other autonomy technology previously invested in by NASA, for example the Smart Executive, can be reused on the ISPP testbed, thus increasing capability without greatly increasing cost.

A longer term goal is to continue research into control of physical systems which must continuously adjust their operation to unforeseen degradation in capability (for example an ISPP unit where Martian dust covers solar panels or slowly clogs air filters) rather than taking a discrete recovery action as Livingstone does. Related issues include reasoning about hybrid discrete/continuous systems, predictive diagnosis and relearning models of the continuous dynamic behavior of the system. This research should contribute to development of ISPP and other robust systems that run at the ragged edge of optimality throughout their lifetimes, neither being overly conservative nor exceeding their remaining degraded capabilities.

Autonomous Rovers

The Remote Agent system, described above and consisting of a planner, a smart execution system and Livingstone, is being adapted for use on the NASA Ames Marsokhod rover as part of an effort to demonstrate increased rover autonomy. That effort is described in [17].

ACKNOWLEDGEMENTS

This paper touches on the work of a great many people too numerous to name here. The Autonomous Systems Group at NASA Ames Research Center consists of about twenty computer science researchers pursuing all manners of autonomy research, much of which was not mentioned here. Members of the JPL New Millennium Program and AI Group contributed to the Remote Agent architecture and to making it work on a flight platform. Advanced development groups at NASA JSC (3T and Bioplex PGT), NASA KSC (ISPP and KATE) and JPL (space based interferometry) have shared their expertise with us and are helping to push the model-based autonomy technologies described here forward.

REFERENCES

Many of the following papers can be found at <http://ic-www.arc.nasa.gov/ic/projects/mba/>

1. R. Zubrin, R. Wagner, 1996. The case for Mars: The plan to settle the Red Planet and why we must. The Free Press.
2. S.J. Hoffman, D.I. Kaplan, Eds., 1997. Human Exploration of Mars: The Reference Mission of the NASA Mars Exploration Study Team. NASA Special Publication 6107. July.
3. A.H. Mishkin, J.C. Morrison, T.T. Nguyen, H.W. Stone, B.K. Cooper, B.H. Wilcox, 1998. Experiences with operations and autonomy of the Mars Pathfinder microrover. Proc. IEEE Aerospace Conf., Snowmass, CO.

4. G.M. Brown, D.E. Bernard, R.D. Rasmussen, 1995. Attitude and articulation control for the Cassini Spacecraft: A fault tolerance overview. Proc. 14th AIAA/IEEE Digital Avionics Systems Conf., Cambridge, MA, Nov.
5. B.C. Williams, P. Nayak, 1996. A Model-based Approach to Reactive Self-Configuring Systems, Proc. AAAI-96.
6. B.C. Williams, B. Millar, 1996. Automated Decomposition of Model-based Learning Problems. Proc. QR-96.
7. N. Muscettola, B. Smith, C. Fry, S. Chien, K. Rajan, G. Rabideau, D. Yan, 1997. Onboard Planning For New Millenium Deep Space One Autonomy, Proc. IEEE Aerospace Conference.
8. B. Pell, E. Gat, R. Keesing, N. Muscettola, B. Smith, 1997. Robust periodic planning and execution for autonomous spacecraft.
9. B. Pell, D. E. Bernard, S. A. Chien, E. Gat, N. Muscettola, P. P. Nayak, M. D. Wagner, B.C. Williams, 1997. An Autonomous Spacecraft Agent Prototype, Proc. 1st Inter. Conf. on Autonomous Agents.
10. D. E. Bernard et al., 1998. Design of the Remote Agent Experiment for Spacecraft Autonomy. Proc. IEEE Aero-98.
11. J. de Kleer, B.C. Williams, 1987. Diagnosing Multiple Faults, Artificial Intelligence, 32(1).
12. J. de Kleer, B.C. Williams, 1989. Diagnosis With Behavioral Modes, Proc. IJCAI-89.
13. J. de Kleer, B.C. Williams, 1991. Artificial Intelligence, 51, Elsevier.
14. S. Weld, J. de Kleer, 1990. Readings in Qualitative Reasoning About Physical Systems, Morgan Kaufmann Publishers, Inc., San Mateo, California.
15. D. Schreckenghost, M. Edeen, R.P. Bonasso, J. Erickson, 1998. Intelligent control of product gas transfer for air revitalization.. Abstract submitted for 28th International Conference on Environmental Systems (ICES), July.
16. R.P. Bonasso, R.J. Firby, E. Gat, D. Kortenkamp, D. Miller, M. Slack, 1997. Experiences with an architecture for intelligent, reactive agents. In Journal of Experimental and Theoretical AI.
17. J. Bresina, G.A. Dorais, K. Golden, D.E. Smith, R. Washington, 1998. Autonomous Rovers for Human Exploration of Mars. Proc. 1st Annual Mars Society Conference. Boulder, CO, August.
18. B.C. Williams, P.P. Nayak, 1996. Immobile Robots: AI in the New Millennium. AI Magazine, Fall.
19. B.C. Williams, P.P. Nayak, 1997. A Reactive Planner for a Model-based Executive. Proc. IJCAI-97.
20. N. Muscettola, 1994. HSTS: Integrating planning and scheduling. In Mark Fox and Monte Zweben, editors, Intelligent Scheduling. Morgan Kaufmann.
21. V. Gupta, R. Jagadeesan, V. Saraswat, 1997. Computing with Continuous Change. Science of Computer Programming.

Mining Automation in the Next Millennium: Engineering a Tele-operated Load Haul Dump Model

Yeen Shien Hwang*¹, Neda Farmer², Jason Hart****

* University of British Columbia, Dep't. of Mining and Mineral Process Engineering,
Vancouver, BC, Canada

¹ Huckleberry Mines Ltd., Houston, B.C., Canada

² Luscar Coal Mine, Hinton, A.B., Canada

** Nautilus International Limited, Burnaby, B.C., Canada

ABSTRACT

A 1:5 scale remote controlled Load Haul Dump (LHD) vehicle has been built to demonstrate the benefits of safety and automation in underground mines. A LHD vehicle is a low profile scooptram that loads broken materials after a blast, hauls the broken rock to a central processing area or conduit where it dumps the material. In most bulk mining operations, LHDs are used inside stopes up to 40 meters high, where tons of rock may fall from above. Controlling an LHD inside the stope from a distance is much safer than drawing broken material while seated on the vehicle. LHDs also pose an accident risk to operators and other bystanders along the travel path of the vehicle. Tele-operated LHDs have higher utilization since a single operator can control one or more LHDs from surface. In such cases, the operating speed of the vehicle can be increased as operator safety is no longer an issue. Hence, higher productivity can be generated. This paper discusses the goals, design and construction of the model and simulates the application of expert system to control the model.

INTRODUCTION

Application of automation in underground mining has grown significantly over the past few decades. In particular, remote-controlled Load Haul Dump (LHD) vehicles are receiving attention because safer working conditions can result together with an increase in productivity. This paper describes development of a 1:5 scale working model of a tele-robotic Load Haul Dump (LHD) vehicle. The model was built to demonstrate the concept of tele-robotics for LHD equipment and to study different control techniques.

The authors who are recent graduates from the University of British Columbia, initiated the project. In September 1997, they approached Nautilus International Ltd of Burnaby, B.C. with the idea to sponsor the project. Nautilus International has extensive experience with mine automation systems and has developed a number of full-scale LHD robotic vehicles at several mines in Australia. The two fourth year students were able to convince senior management at Nautilus International that building such a unit could be useful to demonstrate the conceptual feasibility of automation, as well as, to study alternative control strategies. This project demonstrates how successful collaboration between industry and university students can be achieved. The model was displayed at several events in 1998:

- The UBC Engineering Ball at the Hotel Vancouver in January.
- The Vancouver Branch CIM Student Exhibition Competition at Stanley Park Pavilion in February.
- The Engineering Week Public Demonstration at the Vancouver Public Library in March.
- The CIM 100th Anniversary Annual General Meeting in Montreal in April.

This paper discusses the goals, design and construction of the model and presents a simulation of the application of expert system to control the model autonomously.

MINE AUTOMATION

To survive in today's mineral industry, mining must be both innovative and cost effective to compete with low cost producers in third world countries. It is expected that mining automation technology will yield significant improvements in productivity, efficiency and safety. The benefits of automation are as follows:

- Improved safety for mine workers.
Tele-operation or full automation of LHD can remove the operator from a hazardous environment. For mining operations in underground uranium, coal and asbestos mines, where radioactivity, toxic gas and fibre dust respectively are health concerns for miners, tele-operation can be an ideal solution.
- Increased utility of equipment.
The ability to control various LHDs remotely allows one operator to handle more than one vehicle and to work productively for longer times during a shift.
- Reduced maintenance costs through continuous monitoring.
Fully or semi-autonomous equipment is equipped with many sensors to monitor and control the equipment. Thus, preventive maintenance can be done prior to mechanical failure.
- Increased productivity.
Since the operator is located at a safe distance away from the LHD, the LHD can be operated at higher speed and for longer percentages of each shift.
- Better cash flows from higher throughput.
Higher productivity and utilization ensure higher cash flow for the operation.

The drawbacks of mining automation include:

- Difficulties in economic justification.
It is hard to place economic value on safety, productivity, utilization and maintenance of equipment.
- Acceptance of this new approach by unions and/or workers.
Commissioning of automation can be perceived as a threat to those workers who are ill-equipped to understand the technology and benefits of mine automation.

There are three levels of vehicular automation: line of sight tele-operation, out of sight tele-operation and fully autonomous operation. Line of sight operation allows worker to operate the equipment at a distance away from danger. However, poor visibility may hinder the worker from operating the equipment properly. For out of sight tele-operation, the equipment is equipped with cameras and sound detectors that send images and sound in the vicinity of the LHD back to the operator who may be located on surface. The final level of automation allows the vehicle to drive itself. With this level, the LHD requires virtually no supervision and can navigate itself along known pre-determined paths and can detect and avoid obstacles.

Currently, tele-operation enables underground mines to modify their mining methods to better suit the environment and improve safety and economics. Current examples of such new bulk mining methods are longhole stoping, vertical retreat mining and vertical retreat pillar recovery [1].

LOAD HAUL DUMP (LHD) VEHICLES

A LHD vehicle is a low profile scooptram that loads broken ore after a blast, hauls the material to a central processing area or conduit where the material is dumped. Most underground units are articulated in the centre with a front-mounted scoop which can be raised, lowered and turned to dump its contents. An LHD can be hazardous to underground workers because of dangerous rock-fall conditions and/or narrow haulageways. In most bulk operations, LHDs retrieve broken ore from stopes, that may be up to 40 meters high and then move along a tunnel to dump the ore at an orepass. Operators arriving for their 8 hour shift may often spend up to 45 minutes to get from surface to the machine and another 45 minutes to return to surface by use of a hoist. This "lost" time can be garnered through application of tele-operations.

MODELING A TELE-OPERATED LHD

The LHD model, shown in Figure 1, was designed and built to a scale of 1:5 with dimensions of 2.5m long, 0.5m wide and 0.7m high. Construction of the model began in early October 1997 following on from the engineering and design work conducted in September 1997. All electronic components and instrumentation were installed by Nautilus International in December 1997 and the model was completed in January 1998.

The benefits of the model are as follow:

- Demonstration and Exhibition
The model can be used to present the concept of mine automation.

- **Research and Development**
The model can be used as a research tool to develop navigation, detection, automatic guidance systems and algorithms for operational control of a fleet of LHDs.
- **Training Tool**
New operators can be trained using the model without employing a full-scale unit. Damages incurred to the model during training will be less costly to repair than those of a full-scale vehicle.
- **Study of Ergonomics (Human Factors)**
The study of human factors can improve operational efficiency, health and safety of the operators and the relationship between workers with their equipment and work place.



Fig. 1. The 1:5 scale tele-robotic LHD model developed by UBC and Nautilus International.

ENGINEERING DESIGN

In early September 1997, engineering design began at Nautilus International's office in Burnaby, B.C. With the interactive participation of all members of the design team from both UBC and Nautilus, the design was completed in one month. The model was drawn up in AutoCad as shown in Figure 2.

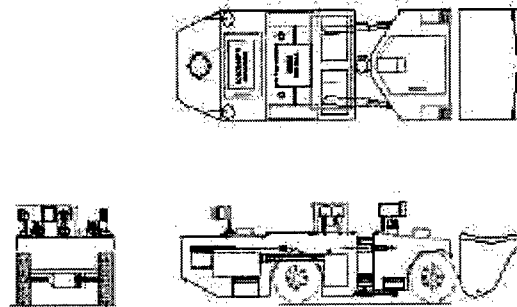


Fig. 2. AutoCAD drawings of the 1:5 scale LHD model.

The model was built from scratch using a second-hand motor and differential retrieved from an electric scooter. We used sheet metal, square metal tubing and other fabricated parts to construct the frame. Unlike a full-scale LHD which is powered by a diesel engine, our model used two 12 volt deep-cycle car batteries.

Consideration of size and weight were crucial during the initial design as we wanted to transport the model to Montreal and other places by plane. Cost was also a key issue. For a full scale LHD, powerful hydraulic actuators control steering, loading and unloading. Due to high cost and complexity of these systems and potential problems such as oil leaks, we substituted mechanical devices for hydraulic actuators.

Chassis

Since the model was to be transported by plane, the chassis was designed as three separate components for quick assembly/disassembly. These components can be seen in Figure 3. The rear and middle components are connected by bolts while the middle and front components are connected by two locking pins. The bucket constituted a fourth component which can be quickly attached to the front of the vehicle chassis.

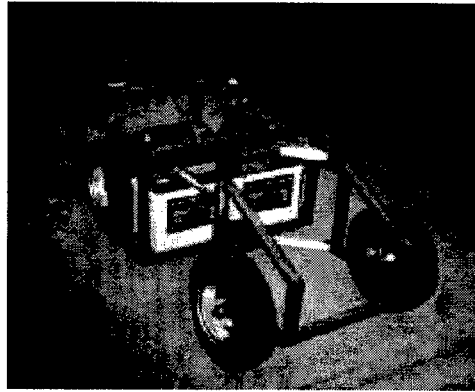


Fig. 3. Chassis of the model.

The chassis was constructed from one-inch square stainless steel tubing. Stainless steel material was chosen since it is strong, rust resistant and easy to weld. The body of the model consists of 14 gauge stainless steel plate. For esthetics, various labels were placed on the body surface.

Bucket

The initial design uses two electrical actuators for controlling the lifting and tilting of the bucket. In this design, three arms were required: two to raise and another to tilt bucket. This design was abandoned due to high cost and other complications. Consequently, another design, as shown in Figure 6, was adopted. Since the model was not designed to load heavy material, the design used a bucket arm to raise the bucket and a home-made actuator to tilt the bucket for dumping. To raise the bucket, the motor drives the chain, which is welded to the bottom of the bucket arm, to lift the bucket. For the tilting operation, a homemade screw-driven actuator was fabricated using a threaded car-jack with an electric motor. The designed stroke for the actuator was 30cm. Electrical brakes were programmed to limit the lift and tilt operations of the bucket.

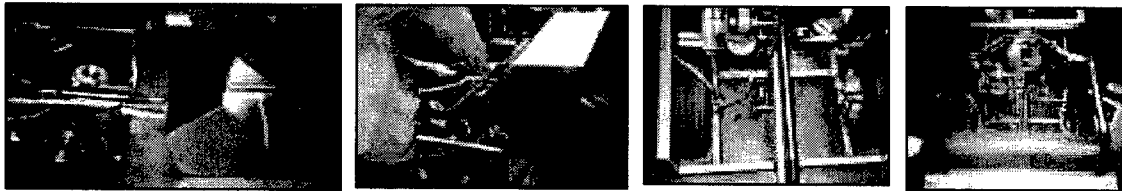


Fig. 6. Bucket design for loading and dumping operation.

Steering Control

The model is articulated similarly to a full-scale unit. The articulated joint is exploited using a mechanical device to provide steering control. A number of alternate designs were examined with consideration given to simplicity, cost and precision. Cable, chain, belt and combinations of such options were tested.

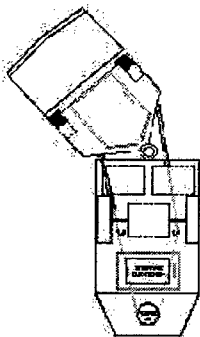


Fig. 4. Articulated LHD model



Fig. 5. Articulated center and mechanical steering limits.

We found that cable and chain were unsuitable for the application because of slipping and skipping. So, it was concluded that the best design option was to use a belt. As shown in Figure 4, the chosen design used an electrical motor, mounted on the front section, to steer the back section using a belt. The front and back section of the model was connected at two points using one-inch pin as shown in Figure 5 below. To ensure that the steering will not exceed 45 degrees from center to side, mechanical obstructions were used to halt the steering to prevent collision.

Wireless Communication

Once the steering system of the model was constructed, Nautilus International installed and wired the electronic equipment. The model is equipped with three cameras, two transmitters, a receiver, four lights and three electronic boxes. Communication between operator and model takes place by wireless communication using Nautilus's portable control unit (PCU) with a distance range from 200 to 400m. This unit can control up to 6 vehicles. Data and video signals are transmitted every 25ms from the model to the PCU. The operator transmits command signals from the PCU back to the model. The video and data signals are transmitted at a frequency of 1300MHz and 500MHz respectively.

Skill is required to drive the model with a joystick. To make it easier to operate, Nautilus International converted the driving system from joystick on the PCU to a PC-based control unit with wheels and pedals, similar to that of an automobile. This control station provides a comfortable workspace that can be located well removed from the unit itself. The guidance system uses video images together with a positioning system overlaid on a map of the underground mine to locate and control the vehicle.

APPLICATION OF AN EXPERT SYSTEM IN THE CONTROL STRATEGY

An expert system was developed to investigate control strategies for the vehicle. The system simulates the autonomous control of the vehicle travelling through a narrow tunnel using Excel. The required sensory information for the control system is present and past wall distance difference of the vehicle to tunnel walls, height and location of obstacles and travelling velocity on each data-transmission cycle received from sensors located on the LHD. The expert system uses fuzzy logic to return the following outputs: turning angle, detection of obstacle, coordinates of the vehicle and speed control.

Since a LHD is articulated at the center, tight steering control is essential. Consequently, fuzzy logic is employed. The best tramming strategy for the vehicle would be to travel along the centerline of the tunnel. Consequently, the difference in the distance between the left and right side of the vehicle to the tunnel walls is used, as shown by equation 1 below:

$$\text{Wall distance difference (} W_d \text{)} = \text{Left wall distance (} W_L \text{)} - \text{Right wall distance (} W_R \text{)} \quad 1.$$

The smaller the wall distance difference, the closer the vehicle is to the left wall. A fuzzy set defining the wall distance difference of the vehicle in the tunnel, as shown in Figure 7, is used to position the vehicle along the centerline of the tunnel. The tunnel is sliced into 5 regions from left to right: Big negative, Small negative, Zero, Small positive and Big positive.

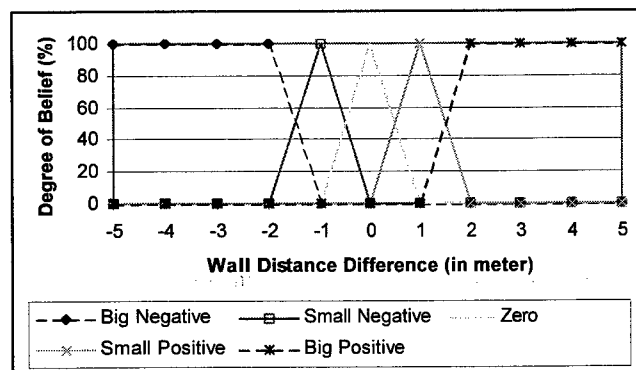


Fig. 7. Fuzzy set definitions for wall distance difference used to determine the turning angle (θ).

At different values of wall distance difference, each fuzzy set has a different degree of belief. The one with the highest degree of belief is dominant in determining the steering angle. As illustrated, fuzzy logic allows the control system to work with uncertainty. A confidence level ranging from 0 to 1, is used to restrict the influence of each fuzzy set to determine the turning angle θ . Only those with degrees of belief greater than the confidence level is applied to the calculation. For this example, a confidence level of 0 has been used.

To minimize drastic turning angles, the system remembers the previous wall distance difference fuzzy sets. Using both previous and present wall distance differences, the system is able to determine the appropriate turning angle to center the model as shown in Table 1.

Table 1.
Determination of turning angle using previous
and present wall distance difference fuzzy sets.

		Past Wall Distance Difference				
		Big Positive	Small Positive	Near Zero	Small Negative	Big Negative
Present Wall Distance Difference	Big Positive	Big Left	Small Left	Zero	Small Left	Zero
	Small Positive	Small Left	Small Left	Small Left	Zero	Small Right
	Near Zero	Small Left	Small Left	Zero	Small Right	Large Right
	Big Negative	Small Left	Zero	Small Right	Large Right	Large Right
	Small Negative	Zero	Small Right	Large Right	Large Right	Large Right

Fuzzy sets were also created to define the turning angle. These were small-right, small-left, zero, big-right and big-left, which are located at supremum positions of 5, -5, 0, 12 and -12 degrees respectively. Weighted-average defuzzification is applied to calculate a discrete turning angle for the steering motor.

In the simulation program, an object of certain height is assigned to a fixed coordinate position in a tunnel. If the object lies in the path of the LHD, the vehicle will detect it within a distance equivalent to 1.5 times its maximum operating speed. If the height of the object is greater than 20 cm, the object becomes an obstacle. Once an obstacle is detected, the control system examines other alternative routes to avoid colliding with the obstacle. If the obstacle is located closer to the right-hand wall, the control system sets the position of the right wall to the location of the obstacle and vice-versa. This forces the vehicle to steer around the obstacle. After passing the obstacle, the location of the wall is returned to its correct position.

If the initial speed of the vehicle is zero, then the system accelerates the LHD up to its maximum speed. The operating speed remains at such unless an obstacle is detected. When an obstacle is detected, the speed of the vehicle is automatically reduced to a recommended speed based on Equation 2:

$$\text{Recommended speed} = \text{ABS} (\Delta W_d / 2 \sin \theta) \quad 2.$$

Also, when an obstacle is detected, the turning angle is magnified by a factor of 5. This allows the vehicle to react faster to avoid a collision. Once the LHD has passed the obstacle, the speed is set back to its maximum value which, in this study, is 4.2m/s.

The above strategies were programmed in Visual Basic and simulated using Excel. Thirty cycles of simulation were conducted in each test, each cycle representing one second of real time. Many scenarios were tested and the expert system was able to navigate the vehicle around any detected obstacle and travel along the centerline of the tunnel or navigate around such obstacles.

CONCLUSION

Mine automation is the future for mining operations to achieve economies of scale. Its benefits are higher productivity, higher equipment utilization, and higher throughput and, last but not least, a safer working environment. This project has successfully demonstrated the usefulness of a fully operational model for use in research and development. As well, the project demonstrates the potential to apply an expert system in navigating an LHD through underground tunnels.

ACKNOWLEDGEMENT

This project was a success due to the dedication and hard work of many participants and helpers. In particular, the authors wish to thank Sean Dessureault, Douglas Bates and Nautilus International. Nautilus International funded this project in collaboration with UBC/MMPE. Also, special thanks to Dr. J. Meech and Dr. M. Scoble for their guidance and support in this project.

REFERENCE

1. Vagenas, N., 1998. Advanced mining technologies: productivity tools for mining in the 21st century. CIM Bulletin, January, p16-19.

Dynamic Reconfiguration of Holonic Lower Level Control

Xiaokun Zhang and Douglas H. Norrie

Division of Manufacturing Engineering, The University of Calgary
2500, University Drive, Calgary, Alberta, T2N-1N4, Canada
Tel: (403) 220-5787 Fax: (403) 282 8406
Email: xkzhang@enme.ucalgary.ca ; norrie@enme.ucalgary.ca
URL: <http://imsg.enme.ucalgary.ca>

ABSTRACT

In this paper, a new approach to dynamic reconfiguration of holon controllers is presented. Based on metamorphic mechanisms for distributed decision-making in agent-based manufacturing systems, the concept of the dynamic virtual cluster is extended to manufacturing process control at the lower levels. Event-driven dynamic clustering of resource control services and cooperative autonomous activities are emphasized in this approach

INTRODUCTION

In recent years, the Holonic Manufacturing System (HMS) has been proposed as an advanced system architecture for intelligent manufacturing systems (IMS). A HMS is composed of different kinds of holon which are autonomous, self-reliant manufacturing-related entities. A holon is an identifiable part of a manufacturing system that has a unique identity, yet is made up of sub-ordinate parts (also holons) and in turn is part of a larger whole (also a holon). A holon consists of an information processing part and often a physical processing part. Autonomy, cooperation, and organizational self-adaptation are considered to be basic characteristics of an HMS.

Holonic architectures and related properties, including autonomy, cooperativeness, and recursivity have been considered by Gou et al. [1], Mathews [2], Brussel et al. [3], and Bussmann [4]. An agent-based view of a holon was suggested in Maturana et al. [5]. A basic concern for an HMS organization is how the resources can be organized dynamically during run-time of the HMS and how the associated controller components can be reconfigured dynamically as well. Intelligent manufacturing is an important application for holonic control processes. Recent research has investigated holonic architectures at the factory or cell level in this area, but relatively little work has been reported for the lower control levels.

In this paper, a new approach to dynamic reconfiguration of holon controllers is presented. Based on metamorphic mechanisms for distributed decision-making in agent-based manufacturing systems [5], the conception of the dynamic virtual clustering is extended to manufacturing process control at the lower levels. Event-driven dynamic clustering of resource control services and cooperative autonomous activities are emphasized in this approach. The paper is organized in two parts. First, the mediator-based dynamic virtual clustering mechanism is presented. Second, the task-driven scheduling and control architecture is introduced and the relevant implementation approach is detailed.

DYNAMIC VIRTUAL CLUSTERING

Dynamic virtual clustering is a dynamic mechanism for organizational reconfiguration of the manufacturing system during run-time. An organization based on virtual clusters of entities can continually be reconfigured in response to changing task requirements. These tasks can include orders, production requests, as well as planning, scheduling, and control. A cluster exists for the duration of the task or sub-task it was created for and is destroyed when the task is completed. Mediators play key roles in the process and manage the clusters. Instead of having pre-established and rigid layers of hierarchical-organized mechanisms, such a mediator-based HMS can use reconfiguration mechanisms to dynamically organize its manufacturing devices. The necessary structures of control are then progressively created during the planning and execution of any production task. In this dynamically changing virtual organization, the

partial control hierarchies are dynamic and transient and the number of control layers for any specific order task are task-oriented and time-dependent.

GT-based Manufacturing Machine Regrouping

The traditional approach to machine layout in a manufacturing system has been predominantly functional (process oriented). Sections of a factory specialize in a particular process or sub-process. Parts requiring more than one process are transported from one section to another until they are completed. The functional layout has a number of disadvantages [6]. Long and uncertain throughput time is a major problem that in turn translates to a high inventory holding cost, untimely product delivery, and increasing losses of sales. Group technology (GT) facilitates an alternative layout of machines in a manufacturing system, which promises reduction of material handling time, queuing time, throughput time, setup time and simplification of tooling [6,7]. GT can be applied to a manufacturing system in two ways: logical or physical. In the logical layout, machines are dedicated to part families but their positions in the manufacturing system are not altered. In the physical machine layout, dedicated machine manufacturing cells containing different machines are created for part families to exploit manufacturing system efficiency [6]. In the intelligent manufacturing system, production order-oriented dynamic grouping of machines promises similar benefit. In this case, deriving and implementing the logical machine layout is the major issue. However, during run-time, the changes in the dynamic groups of the machines require that the control levels have appropriate flexibility in function and infrastructure.

Control System Reconfiguration

Intelligent manufacturing systems require dynamic reconfigurability at all levels of all system components. This encompasses online changes at the hardware, network communications, system software and application software levels. Online change requires recognition of such changes immediately and acting accordingly. Increasing the autonomy of individual components within a manufacturing system diminishes the number of control levels compared to a conventional hierarchical configuration [8]. Dynamic grouping requires to have the components autonomy and can thus be associated with fewer levels of control. As will be seen later, the controller configuration involves a production-task-oriented controller cluster, which facilitates collaboration and reduces communication delays during task scheduling and control cycles. This approach facilitates control fault-tolerance especially within a hardware-redundant environment.

Mediator Clusters

A basic HMS architecture can be based on four holon types: Product Holon (PH), Product Model Holon (PMH), Resource Holon (RH), and Mediator Holon (MH). A Product Holon holds information about the process status of product components during manufacturing, time constraint variables, quality status, and decision knowledge relating to the order request. A Product Holon is a dual of a physical "component" and an information "component". The physical component of the Product Holon develops from its initial status (raw materials or unfinished product) to an intermediate product, and then to the finished one, i.e. the end product. A Product Model Holon holds up-to-date engineering information relating to the product life cycle (configuration, design, process plans, bills of materials, quality assurance procedures, etc.). A Resource Holon contains physical and information components. The physical part contains a production resource of the manufacturing system (machine, conveyor, pallet, tool, raw material, and end product, or accessories for assembling etc.), together with controller components. The information part contains planning and scheduling components (see later section for details).

In contrast with some other HMS models [1,3,4,9], the concept of the Mediator Holon is emphasized in the approach described in this paper. A Mediator Holon serves as an intelligent logical interconnection to link and manage orders, product data, and specific manufacturing resources dynamically. The Mediator Holon can collaborate with other holons to search for and coordinate resource, product data, and related production tasks. A Mediator Holon is itself a holarchy. A Mediator Holon can create a Dynamic Mediator Holon (DMH) for a new task such as a new order request or sub-order task request. The Dynamic Mediator Holon then has the responsibility for the assigned task. When the task is completed, the DMH is destroyed or terminates for reuse. DMHs identify order-related resource clusters (i.e. machine group) and manage task decomposition associated with their clusters.

Product Holons and Mediators

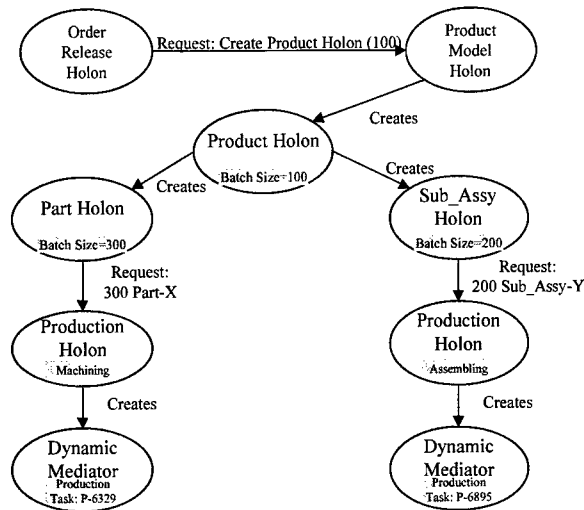


Fig.1. Initial activity sequence following order release

Fig. 1 shows the initial activity sequence following the release to production of an order for 100 of a particular product. This simple example considers the product to be composed of 3 identical parts (to be machined) and 2 identical sub-assemblies (each to be assembled). As seen in Fig.1, following the creation of the appropriate Product Holon, there are created the relevant Part and Sub-Assembly Holons. The requests for manufacturing made by these latter holons to appropriate Production Holons (which function as high-level Production Managers for a manufacturing shop-floor plan or part dispatch) result in the creation of Dynamic Mediators for the machining and assembly tasks. Subsequently, each Production Holon coordinates inspection or assembly of the parts or sub-assemblies according to the production sequence prescribed by the Production Model Holon (from its stored information). More complex situations occur, when products having many components requiring different types of production processes are involved.

Logical and Physical Machine Clusters

After GT-based physical and logical machine groups are derived, the necessary control structures are created and configured using control components cloned from template libraries by a DMH. The machine groups, their associated and configured controllers, then form a temporary manufacturing community, termed a virtual cluster holon (VCH) as shown in Fig. 2. The VCH exists for the duration of the relevant job processing and is destroyed when these production processes are completed. The physical component of a VCH is composed of order-related parts, raw materials or sub-products for assembly, manufacturing machines and tools, and associated controller hardware. Within these manufacturing environments, parts grow from their initial state to an intermediate product and then to the finished one. The information component of a VCH is composed of cluster controller software-components, the associated DMH, and intermediate information on the order and the related product. Each cluster controller is further composed of multi-layer control functions that execute job collaboration, control application generation and controller dynamic reconfiguration, process execution, and process monitoring, etc.).

TASK-DRIVEN SCHEDULING AND CONTROL

Dynamic Virtual Cluster and Controller Cluster

The life cycle of a dynamic virtual cluster holon has four stages: Resource grouping; control components creation; execution processing; and termination/destruction. The Dynamic Mediator Holon is involved in the stages 1 and 2. The first cluster that is created is the schedule-control cluster shown in Fig 2 & 3. Once the cluster is created, it continues, due to its autonomy. A cluster can be also considered to be a holonic grouping. The Controller Cluster next created is composed of three holonic parts: Collaboration Controller (CC), Execution Controller (EC), and Control Execution (CE) holon. One CE holon can be associated with more than one physical controller (execution platform such as real-time operation system and its hardware support devices) and appears like a distributed-node transparent-resource platform for execution of cluster control tasks at the resource level. In the prototype system under development, the CC, EC and CE holons collaborate to control and execute the distributed tasks or applications on a new type of distributed real-time operating system recently implemented [10]. The distributed tasks or applications are represented using the Function Block (FB)-1499 specification, which is a draft standard described by the IEC for distributed industrial-process measurement and control systems.

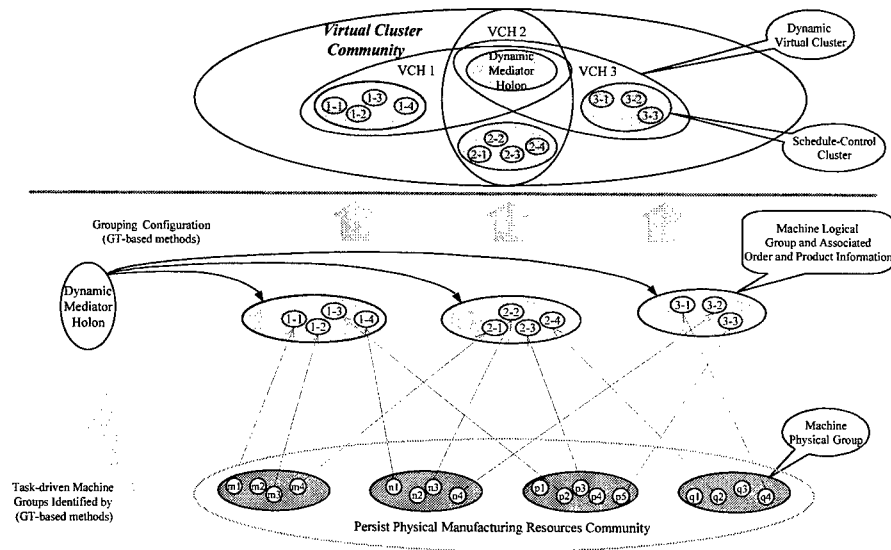


Fig. 2 Machine physical and logical grouping with dynamic virtual clustering

Resource Scheduling and Control

Fig. 3 illustrates resource scheduling and control. For simplicity, only production resources such as CNC Machining or CNC Turning Centers are considered (for the more complex case where transport resources need to be coordinated with production resources, see Ref. [11]). The following describes the entities (holons) shown in the Fig. 3 and what their responsibilities are. The Collaboration Controller (CC) holon receives production requests from the Dynamic Mediator of its dynamic virtual cluster (see Fig. 2) and reports back as needed on the status of these tasks and any re-assignment needed. The CC is responsible for building and maintaining the "joint schedule" for the resources under its control. This joint schedule can be thought of as a Gantt chart for the schedules of these machines (a sliding window moving forward in time and covering entire a pre-determined period or pre-determined number of jobs). Each resource has a resource planner for which a (partial) Clone is created and assigned to the Schedule-Control Cluster. Each resource will have a clone in each such cluster it is involved with. The Collaboration Controller sends production requests (includes job number; part or assembly ID; quality; due date; etc.) to the Resource clones. Each clone then negotiates a "provisional schedule" with its Resource Planner. This pre-assumed schedule is then checked with the Resource Scheduler, by the Resource Planner, and is returned as a "committed schedule" or a "modified provisional schedule" (best fit). In the latter case, this new provisional

schedule is passed back to the Resource clone for re-negotiation with the Collaboration Controller (which checks this against the production request and its joint schedule). By such mechanisms, the schedule for each resource is progressively determined and communicated to the Collaboration Controller which updates its present schedule accordingly and takes appropriate action (to inform the Dynamic Mediator of out-of-date schedules; to send execution requests to lower level controller, etc.).

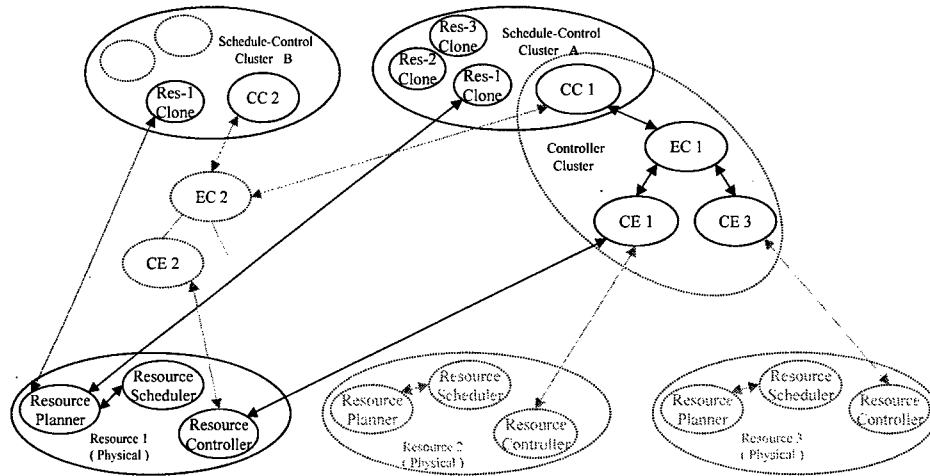


Fig. 3 Resource scheduling and control

We now consider what happens when the Collaboration Controller issues an execution request to an Execution Control Holon. This downloaded information is essentially a request to start a named job on a specified resource at a given time. This job-execution high-level "control code" is an "application" for the EC to arrange execution. It does this by preparing a function-block representation of the application procedure, using application prototype (templates) and function-blocks (both composite and basic) from libraries. It then arranges for compilation into low-level execution control code and distributes the application modules of the code among the related Control Execution Holons. Each Control Execution Holon in turn arranges for these execution modules to be distributed in suitable form to the Resource Controllers of the appropriate resources. A Resource Controller is a software Holon paired with a Physical Controller that executes this resource-level application control code.

If there are problems at the machine level (e.g. setup delayed due to a broken tool; machine shut down due to overheating or excessive vibration), this information goes up to the CE which may pass it on to the EC from which it may go further. Each level of control has responsibility for certain levels of remedial action. The CE, for example, has responsibilities for basic monitoring and alerts, for fault-recovery and remedial action. When delays are likely to require job rescheduling, this is handled via a CC and the scheduling procedures described previously.

It is of interest to consider what to happen if say Resource-1 has a clone also in another Schedule-Control Cluster (e.g. Schedule-Control Cluster B shown in Fig. 3). Suppose this clone negotiates with Collaboration Controller (CC 2) for scheduling a job on Resource-1. The execution of this job will then be controlled by CC2 through a new dynamic-created controller Cluster (of CC2, EC2, and CE2...). For the execution of this job, Resource-1 will be under the control of a newly assigned CE (say CE2). When it is completed, the next job held in the "master schedule" of its Resource Scheduler will be dealt with next. If the clone in the original Schedule-Controller Cluster had negotiated this job, it will be under the control of CC1, which may use the original Controller Cluster (or create a new or different one if that accords with current needs).

Thus, we see that scheduling and control are dynamic interlinked activities involving clones, mediators (each CC is of type mediator), dynamically created clusters, and different types of controllers.

Task-driven Activity or Process Sequences

The section above has illustrated the activity (process) sequence using an example. The following gives an overview of sequences for the more general case. Whenever order, resource, or product data changes, this initiates consequential activity within the system. Consider the situation following an order request event:

Step1: An Order Release holon sends an order request to an appropriate high level Mediator holon, and the Mediator holon creates a new Dynamic Mediator holon.

Step2: The Dynamic Mediator holon checks the capabilities of registered resources in its resource library or elsewhere, then executes manufacturing-resource grouping algorithms to form primary resource groups. At the same time, the order is decomposed into component production tasks or sub-tasks. Obtaining this information assists the subsequent order and production task negotiation process. It is a temporary centralized strategy that enhances the effects of the contract net-based order dispatch mechanism among order release holons and interested manufacturing resources.

Step3: The Schedule-Control clusters are next generated and populated with the appropriate resource clones. The Schedule-Control cluster serves as an information part of the Dynamic Virtual Cluster holon and the grouped manufacturing resources serve as a physical part.

Step4: The control components are next created for clustered resources, i.e. creation of the Collaboration Control (CC) holon, Execution Control (EC) holons and configuration with persisting Control Execution (CE) holons to form a Controller Cluster holon.

Step5: The Dynamic Virtual Cluster holons begin negotiation with relevant holons using the information from the Mediator's grouping solution for this order. This negotiation process is based on contract-net mechanisms.

Step6: When a Dynamic Virtual Cluster has negotiated a job and knows the associated resources, its Schedule-Controller Cluster supervises the scheduling process described previously. Then the high-level control tasks are generated by its CC.

Step7: The Execution Control holon generates the tasks or application code for the specific controller platform. These control tasks are represented using the FB-1499 specification.

Step8: The EC holon downloads the application code to the CE platform and initiates execution.

MULTI-AGENT BASED ARCHITECTURE FOR CONTROLLER CLUSTER

Fig. 4 shows the architecture for the Controller Cluster based on a new concept of a CC/EC/CE multi-layer intelligent controller. The information part of each holon (Controller Cluster relevant holons) is based on multi-agent architecture and some of the component code is based on mobile agent mechanisms. Some other implementation details are also shown in the fig. Future development and implementation based on the new intelligent controller concept is in progress.

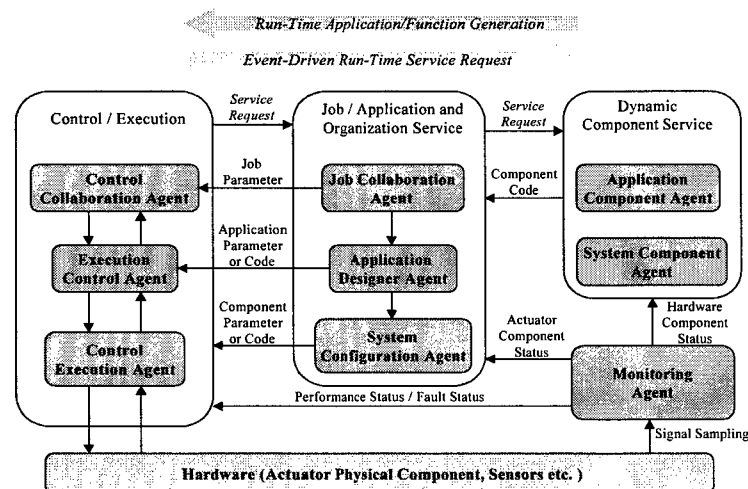


Fig. 4 Multi-agent based architecture for dynamic Controller Cluster

CONCLUSION

Dynamic architectural mechanisms and intelligent reconfiguration are important issues for holonic low-level control. This paper outlines an approach for mediator-based holonic control at both production and control levels. Task-driven scheduling and control is presented at three levels: knowledge-based manufacturing resources clustering; dynamic system reconfiguration; execution. A multi-agent based Controller Cluster architecture with some implementation details is described. The implementation of a prototype system based on the architectural concepts and implementation approaches presented is in progress.

REFERENCES

1. L. Gou, P. B. Luh, and Y. Kyoya, 1998. Holonic manufacturing scheduling: architecture, cooperation mechanism, and implementation. *Computer in Industrial* 37, 213-231.
2. J. Mathews, 1998. Organization foundations of intelligent manufacturing systems— The Holonic Viewpoint, *Computer Integrated Manufacturing Systems*, 8(4), 237-243.
3. H. V. Brussel, J. Wyns, P. Valckenaers, L. Bongaerts, and P. Peeters, 1998. Reference architecture for holonic manufacturing systems: PROSA, *Computer in Industrial* 37, 255-274.
4. S. Bussmann, 1998. An agent-oriented architecture for holonic manufacturing control, In *Proceedings of the 1st Int. Workshop on Intelligent Manufacturing Systems*, EPFL, Lausanne, Switzerland.
5. F. P. Maturana, D. H. Norrie, 1997. Distributed decision-making using the contract net within a mediator architecture, *Decision Support Systems* 20, 53-64.
6. Kusiak, W.J. Boe, and C. Cheng, 1993. Designing cellular manufacturing systems: branch-and-bound and A* approaches, *IIE Transactions*, 25(4), 46-56.
7. Kusiak, W.J., 1990. *Intelligent manufacturing systems*, Prentice Hall, Englewood Cliffs, N.J.
8. A. R. Chaturvedi, R. Gulati, G. Koehler, 1994. *Computational Ecology of Manufacturing Systems*, AAAI-94 Workshop Program Reasoning about the Shop Floor (SIGMAN), Seattle, Washington.
9. S. Kirn, 1996. Organizational intelligent and distributed artificial intelligent, In G.M.P. O'Hare and N. R. Jennings (Eds.), *Foundations of distributed artificial intelligence*, John Wiley & Sons, NY, 505-526.
10. Xiaokun Zhang, Sivaram Balasubramanian, and Douglas H. Norrie, 1999. An Intelligent Controller Implementation for Holonic Systems: DCOS-1 Architecture, Submitted to 14th IEEE International Symposium on Intelligent Control, September 15-17, Cambridge, Massachusetts, USA.
11. F. P. Maturana, 1997. *MetaMorph: an adaptive multi-agent architecture for advanced manufacturing systems*, Ph.D. dissertation, The University of Calgary.

Intelligent Process Monitoring for Paper Machines

Janos L. Grantner*, Peter E. Parker**, George A. Fodor***

*Department of Electrical and Computer Engineering
Western Michigan University, Kalamazoo, MI 49008-5066, USA

**Department of Paper and Printing Science and Engineering
Western Michigan University, Kalamazoo, MI 49008-5066, USA

***ABB Automation Products AB, S-721 67 Vasteras, Sweden

ABSTRACT

Paper machine control is a complex control environment. It consists of a large number of smaller control units, typically Programmable Logic Controllers (PLCs), that are integrated into an overall control architecture. Basis weight and formation control are typical of such systems, with slice screws, basis weight valves, consistency controllers, and pump speed controllers all interacting to produce a uniform sheet in the machine and cross machine directions. Errors in measurements due to sensor malfunctioning, or process states outside the basic assumptions for the control action lead to unwanted and/or poor supervisory response. Manual control often becomes the only way to return the machine to the desired state. In this paper, we report on the initial results of a research to apply the theories of Ontological Control and the Hybrid Fuzzy-Boolean Finite State Machine (HFB-FSM) to paper machine basis weight control. The objective of the research is to provide an automated error detection and recovery method when control encounters an unexpected change in the process environment.

INTRODUCTION

In many complex real-time industrial applications such as pulp and paper, the control system is made up of a large number of smaller control units (e.g., Programmable Logic Controllers-PLCs) integrated into an overall control architecture. When controllers act in a sequential fashion, the output of one controller may be among the input signals of another controller. It is often recognized that these types of systems are complex due to the size, and the number of the possible state combinations in the total state space. However, when controllers of different makes and heterogeneous types are connected together, even the knowledge about the total state set may not be sufficient for a correct supervision. There are always assumptions, often undocumented, about the conditions under which a controller algorithm can be used such that the intended control goals will be reached [1]. Basis weight and formation control are typical of such systems, with slice screws, basis weight valves, consistency controllers, and pump speed controllers all interacting to produce a uniform sheet in the machine and cross machine directions.

Industrial-strength complex control systems are required to act consistently relative to the initial goals when meeting unexpected situations in their environment. The capacity of a system to identify and recover from an error after meeting an unexpected situation is regarded in industry as a very important property. The research reported here proposes a solution to the problem of extending the safety and recovery capacity of complex control systems by introducing a new type of execution monitoring. The solution employs the theories of the Hybrid Fuzzy-Boolean Finite State Machine (HFB-FSM) [2], and Ontological Control [1]. The main points of the approach are as follows:

- It has been shown in the theory of ontological control that problematic control situations at the reactive level (referred to as *state de-synchronization*) can be formally represented and classified. When the state set is well determined, a recovery operation is possible within certain constraints. The constraints are given in terms of an event-driven dynamic linguistic model implemented by the HFB-FSM. These boundaries can be used to specify the recovery capacity of a control system with a given set of actuator and sensor equipment.
- It has been shown that the recovery operation cannot be performed using the state of the reactive-level controller that needs to be recovered. Thus the method exploits a connected supervisory controller for the recovery operation.

- The fuzzy state boundaries for the HFB-FSM are devised using a continuous model of the system.

When an unexpected change occurs in the environment of a reactive-level controller, a supervisory controller monitoring the reactive-level controller can detect that by using the theory of ontological control. The execution-monitoring unit will then invoke the fuzzy specification of the discrete states that are involved in the erroneous situation. That includes a set of fuzzy states of the HFB-FSM, and an algorithm for triggering transients of fuzzy states. A HFB-FSM can model a hybrid system of continuous and two-valued signals when the status of the system can be viewed as somewhere in-between the scope of the discrete states of a PLC control program. The execution monitoring unit will use the information on the next fuzzy state to determine if a recovery is possible, that is, the HFB-FSM represents the specification of the bounds for a recovery. If the HFB-FSM enters a suitable fuzzy state, it returns the particular control action that will achieve the recovery of the reactive controller. The following two sections present the main theoretical tools involved, and then the method is illustrated by an example to monitor and control a paper machine.

HYBRID FUZZY-BOOLEAN FINITE STATE MACHINE

The HFB-FSM can be implemented by a Boolean automaton based upon two-valued logic. It is defined by the formulas (1), where X_F and Z_F stand for a finite set of fuzzy inputs and outputs, respectively, W_B and U_B stand for a finite set of two-valued logic inputs and outputs, respectively. Defuzzified outputs are denoted by z_c , R^* is a composite linguistic model (3), and O is the operator of composition. Each crisp state of the HFB-FSM is characterized by an overall linguistic model R_s , or by a set of linguistic sub-models in the case of multiple-input-single-output (MISO), and multiple-input-multiple-output (MIMO) systems.

$$\begin{aligned}
 Z_F &= X_F \circ R^* \\
 R^* &= G(R_s) \\
 z_c &= DF(Z_F) \\
 U_B &= f_u(y_B) \\
 X_B &= B(X_F) \\
 Z_B &= B(Z_F) \\
 Y_B &= f_y(X_B, W_B, Z_B, y_B)
 \end{aligned} \tag{1}$$

A fuzzy state is defined by a crisp (Boolean) state and a state membership function

$$S_{F_k} : S_k, g_{S_k} \tag{2}$$

where S_{F_k} stands for fuzzy state k , S_k represents crisp state k , and g_{S_k} is the state membership function associated with S_k . G stands for the matrix of state membership functions, X_B , Z_B , Y_B , and y_B are two-valued Boolean input, output and state variables, respectively. B stands for a Fuzzy-to-Boolean transformation algorithm to map a change in the status of a fuzzy variable into state changes of a finite set of corresponding Boolean variables. The z_c crisp values of the fuzzy outputs are obtained by evaluating a defuzzification strategy, DF . On the basis of the concept of a fuzzy state, the FSM stays in a number of crisp states simultaneously, to a certain degree in each. One of these states is referred to as a dominant state for which the state membership function is a 1 (full membership). For each fuzzy state of the HFB FSM model, a R_i^* composite linguistic model is created from the finite set of R_{S_i} overall linguistic models ($i=1, \dots, p$). Let the HFB-FSM be in fuzzy state S_{F_k} , then

$$\begin{aligned}
 R_k^* &= \max[\min(\beta_1^k, R_{S_1}), \min(\beta_2^k, R_{S_2}), \dots, \\
 &\dots, \min(\beta_k^k, R_{S_k}), \dots, \min(\beta_p^k, R_{S_p})]
 \end{aligned} \tag{3}$$

where $\beta_1^k, \beta_2^k, \dots, \beta_p^k$ stand for the degrees of state membership function g_{S_k} , and $R_{S_1}, R_{S_2}, \dots, R_{S_p}$ are the overall rules in crisp states S_1, S_2, \dots, S_p , respectively. With (3), a SISO system is assumed. In adaptive

systems R_k^* is not stored in memory, it is dynamically created by computing (3), instead. By modifying the β degrees of the state membership functions on-line, new R^* composite linguistic models can be created under real-time conditions. The transition between active composite linguistic models is determined by the state transients of the HFB-FSM.

The state transients of the HFB-FSM are specified by means of a sequence of changes in the states of the fuzzy inputs and outputs, as well as of the two-valued inputs. The changes in the states of the fuzzy inputs and outputs are mapped into the corresponding sequence of changes of Boolean input and output variable sets, respectively, using the B algorithm [3]. In this domain, those changes are joined by the state changes of the two-valued inputs. On the basis of this combined Boolean input/output sequence specification the crisp automaton section of the HFB-FSM will then be synthesized. Hence, the HFB-FSM model allows the integration of fuzzy and two-valued logic specifications to describe a system's behavior. The integrated treatment of fuzzy (continuous) and two-valued signals is of great importance for designing complex systems. The theory of the HFB-FSM will be used in the sequel to devise a proper control action when the web breaks in the paper machine.

ONTOLOGICAL CONTROL

When a control system such as an autonomous agent is designed, the modeling assumptions for its control algorithm are inherently extended by additional assumptions about the complex environment in which the control will take place. These assumptions are not represented by formal means, hence, an agent cannot verify whether they are true. Ontological control investigates the case when these assumptions are violated, situations in which a controller acts under *violations of the ontological assumptions*. The architecture of an Ontological Controller (OC) capable to detect violations of ontological assumptions (VOA) in the context of a de-synchronization from the execution of a goal path was shown in [4]. However, the OC is unable to recover from a VOA by itself [4]. The concept of a state is defined in the OC by (4):

$$S_i = (y_i, u_{ij}) \quad (i, j=1, \dots, n) \quad 4.$$

where y_i stands for a two-valued Boolean formula (referred to as a *plant formula*) showing the condition that is true at a given time in the controlled plant. Each relevant plant situation has a corresponding plant formula in some state. A control action denoted as u_{ij} is executed when y_i is true. The expected outcome of this action is that the plant changes such that at the next time instance y_j will be true. However, if some external action (disturbance) occurs, the expected change in the plant does not take place but a new plant state, y_k , will materialize instead. The disturbance is considered as an external action and, if known in advance, is denoted as $u_{i,k}^{ext}$ where the subscripts refer to the respective two plant formulas before and after the external action. The new state can be denoted as some $S_k = (y_k, u_{k,l})$. The control will then proceed in a succession of states. The states that can materialize from an arbitrary state S_i by external actions (disturbances) are referred to as *collateral states* to S_i and the set of such states is denoted as $K(S_i)$.

It has been shown in that a VOA manifests always as a state transition from S_i to a state in $K(S_i)$ where S_j is the consecutive (next expected) state to S_i . This type of transition is referred to as an ontological de-synchronization. The recovery operation can be performed if two conditions are met:

- (i) The cause for the loop is recognized. That is accomplished by the theory of ontological control.
- (ii) Overload limits can be specified for the actuators. That is, the boundaries of a state can be extended.

That is achieved by using the theory of the HFB-FSM.

The recovery from an erroneous control cycle requires the determination of a state that corresponds to the current plant situation, and it has a control action that can "break" the cycle. For fuzzy controllers, there can be found recovery solutions by adding extra rules to the rule base, that is, reacting to the fact that the model of the plant has changed [5]. However, for controllers such as PLCs that have discrete states, the problem of recovery becomes more difficult since there is no state with acceptable properties in the state space of the controller that can materialize at a VOA. Thus the recovery approach suggested in this paper relies on techniques that can accommodate fuzzy states and produce both continuous and two-valued outputs.

RECOVERY APPROACH WHEN THE WEB BREAKS

Modern paper machines produce continuous webs 6 to 9 meters wide at speeds of 1500 - 2000 m/min. Dry web weight (basis weight) ranges from that of light tissue papers at 8 - 10 gm/m² to very heavy board product at 90 - 100 gm/m². Unfortunately no sensors currently exist to monitor on-line, in real time, product attributes that are of prime importance to final customers. These properties, such as machine direction (MD) vs. cross machine direction (CD) strength (which is controlled by the relative orientation of fibers as they are deposited on the wire), formation, or MD and CD basis weight variation (which is primarily controlled by the slurry concentration and turbulence), and various surface attributes can only be measured off-line and result in a significant delay between production and measurement. On line sensors are in development that will permit some of these properties to be measured. However, none is a standard component of a paper machine control system at the current time. The papermaker, then, must rely upon indirect measures to control the machine. The most common method is to measure the moisture and basis weight of the sheet just prior to the reel. Information from these measurements is then used to control drying (for moisture), slurry consistency (for basis weight), and, to a much lesser extent, CD moisture and basis weight profiles.

Figure 1 is a simplified block diagram of a paper machine that can be used to illustrate the control problem. Fiber is fed to the process from some processing step that makes a slurry that is normally around 4 - 6% consistency. After several cleaning and fiber treatment steps, this slurry is delivered to the main feed, or machine chest at about 3% consistency. Stock from this chest is then fed to the fan pump via the basis weight valve. The main feed to the fan pump is dilution water from the white water chests that hold water that is drained from the machine. The machine chest stock is injected into this dilution water at the fan pump suction, with the fan pump acting as an in-line mixer. The stock is now < 1% consistency and is sent to the headbox via some final cleaning equipment. The dilute stock is spread on the wire at the headbox. The primary function of the headbox is to uniformly distribute the wire, using turbulence to break up any flocs. Water drains from the stock through the wire and is returned to the white water chests. Only 50 - 80% of the fibers are retained on the wire, so this water contains 50 - 20% of the fiber originally delivered to headbox. Thus, the recirculating fibers can be equivalent in mass flow rate to the fiber delivered from the machine chest. Simple and vacuum assisted drainage can only remove a limited amount of water, typically giving a web of about 20% solids (or 80% water). This web is then fed to a set of presses which remove more water, with the sheet leaving the press at 40 - 50% solids. The press water is typically filtered and used on various showers; some is recirculated for stock dilution. The semi-dry web is then fed to a series of dryers where the final water is removed by evaporation. After the dryers, the web finally passes under the scanner heads where basis weight and moisture are measured. It is then reeled.

Delay times between action at the basis weight valve and results at the sensor can be relatively long. The distance from the headbox to the scanner may be as much as 100 m for a machine producing heavy weight paperboard. With a speed of 500 m a minute or so for this heavy weight, the delay time is around 12 - 15 seconds. In addition, the scanner travels back and forth across the web at about 0.3 m/s, so for a 6 m wide machine, the time to obtain a signal (1 complete cycle) is about 40 seconds. Thus, total delay between the time a control action is initiated and the response can be measured may be upwards of 1 minute.

When all runs well, the machines, in spite of the long delays, run well and produce stable product. Occasionally, however, the web breaks and recovery procedures are necessary. When the web breaks several things happen. First, the basis weight sensor loses its input, so it can no longer control the basis weight valve. Second, unless the situation is very serious, the web is produced on the wire but is directed to the broke pit prior to the press. In the broke pit, water from the white water chest is added and the stock is diluted and then sent to the machine chest. Consistency in the machine chest is then upset as it is very unlikely that the stock coming from the broke chest is the same consistency as that in the machine chest. (Machine control attempts to do this, but the dynamics of the process prevent excellent control.) With consistency to the machine chest upset, the mass flow of fiber to the headbox changes. Breaks may last 5 - 10 minutes (and sometimes a lot longer) which is sufficient time for the basis weight control to drift rather significantly from its target. When the web is re-established on the reel, the sensor then resumes control and attempts to bring all measurements back in line. This may take several minutes and the production made in this period is off quality, hence, may need to be rejected.

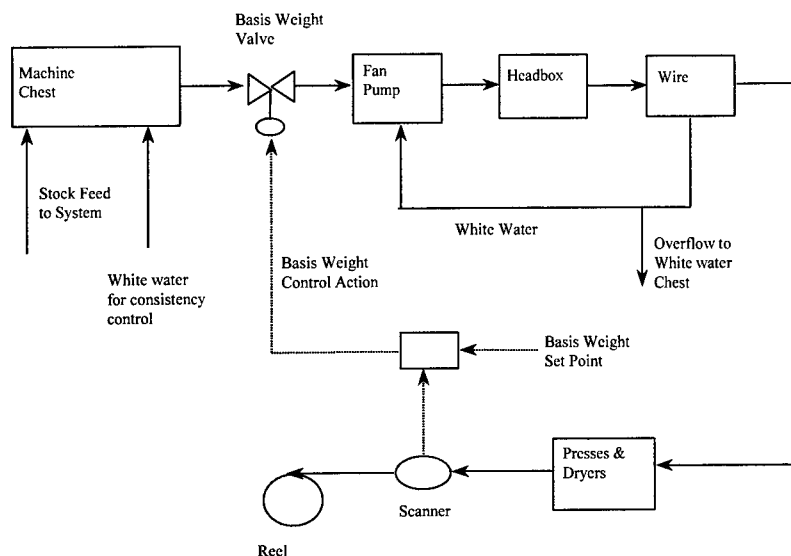


Figure 1. Block Diagram for Basis Weight Control

Much information exists about the state of the system prior to and during the break. Tank levels, flow rates, and various stock consistencies are available. The goal of the recovery operation is to ensure that the consistency and flow rate to the headbox will not drift during a break. Thus, when the web is re-established, the product would be on specification.

Figure 2 illustrates the stock flow during web break. In the normal case, the measurement taken from the scanning sensor is used to control the flow of stock to the fan pump via the basis weight valve. Consistency in the machine chest is kept constant via another control loop that is independent of basis weight. When the web breaks, the sheet goes to the broke pit where it is diluted with white water. Consistency control is problematic, as it is difficult to measure when the sheet is being broken up. Normally, the flow of white water is ratioed to the machine speed and basis weight. The stock is then sent to the machine chest. Since its consistency is probably not the same as the machine chest, the machine chest consistency is thus upset and stock flow to the paper machine changes.

In order to maintain the quality of the product, the following intelligent control approach is proposed: the supervisory controller continuously supplies the HFB-FSM with fuzzified sensory data, hence the fuzzy automaton can monitor the state of the paper making process via appropriate fuzzy state transients. During normal regime, the paper machine is run under the established, classical control algorithm. When the measurements begin to drift due to the web break, the HFB-FSM will move to a corresponding fuzzy state. At this point, using a two-valued output of the HFB-FSM, the supervisory controller shuts off the classical control algorithm and begins to use the continuous (defuzzified) and two-valued outputs of the HFB-FSM to control the paper machine. The required actuator outputs to keep the mass flow of stock ($= \text{flow rate} * \text{consistency} * \text{density}$) constant will be inferred using a knowledge base in that state of the HFB-FSM. (Density is assumed to be constant for this problem.) In other words, a fuzzy model of the paper machine, that has been developed off-line on the basis of expert knowledge and available measured data will be used along with actual data during the break to maintain the consistency of the stock. After the web has been reestablished and the specs are within the required bounds, the supervisory controller will switch back to the classical controller and the HFB-FSM will monitor the status of the process in stand-by.

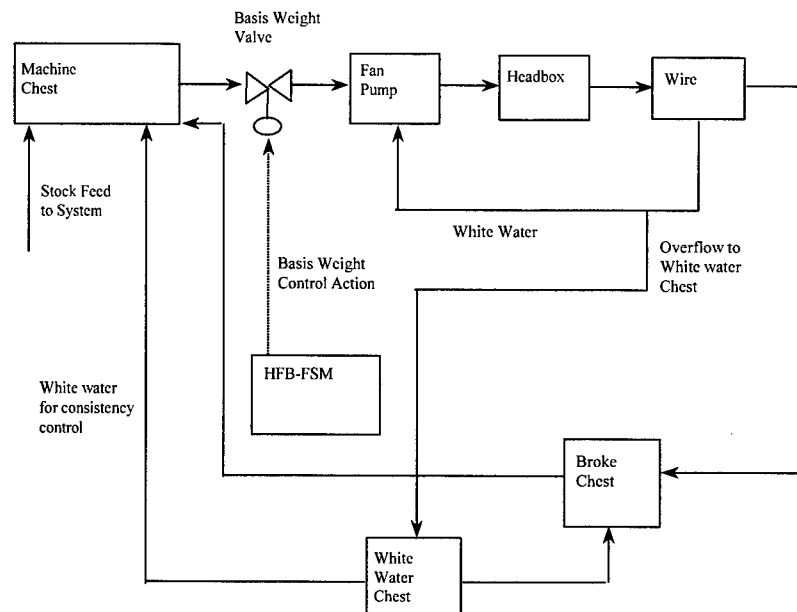


Figure 2. Stock Flow During Break

CONCLUSION AND FURTHER RESEARCH

An intelligent supervisory control approach using a fuzzy automaton was proposed to monitor the status of a paper machine and to maintain the consistency of the stock when the web breaks. The method will be tested in a two step process using Western Michigan University's pilot paper machine. A dynamic model of this machine is under construction using the CADSIM Plus [6] simulation package. A fuzzy automaton model of the control algorithm will also be created using the HFB-FSM Simulator that has been developed at the ECE Department. After the plant model is verified using pilot machine data, breaks will be simulated and the process will be controlled by the HFB-FSM model. Once the HFB-FSM is proven successful on the simulated paper machine, it will be applied to the pilot machine. Being a pilot paper machine, breaks of any duration can be studied and extensive data can be obtained to both verify and tune the HFB-FSM controller.

REFERENCES

1. G.A. Fodor, 1997. *Ontologically Controlled Autonomous Systems: Principles, Operations and Architecture*, Kluwer Academic Publishers, Boston/Dordrecht/London.
2. J.L. Grantner, G. Fodor, D. Driankov, 1998. Hybrid Fuzzy-Boolean Automata for Ontological Controllers, Proc. World Congress for Computational Intelligence, WCCI98, Anchorage, Alaska, USA, May, 1998.
3. J.L. Grantner, 1994. *Design of Event-Driven Real-Time Linguistic Models Based on Fuzzy Logic Finite State Machines for High-Speed Intelligent Fuzzy Logic Controllers*. Dissertation for the Degree Candidate of Technical Science, Hungarian Academy of Sciences, Hungary.
4. G.A. Fodor, 1995. *Ontological Control: Description, Identification and Recovery from Problematic Control Situations*. Ph.D. Dissertation, Dept. of Computer Science, University of Linköping, Sweden.
5. D. Driankov, G. Fodor, 1996. Fuzzy control under violations of ontological assumptions, invited plenary talk. in the Proceedings of the FLAMOC'96 Conference, Sydney, Australia, 109-115.
6. CADSIM Plus, Aurel Systems, Burnaby, BC, Canada

An Integration Design Approach in PID Controller

Jen-Yang Chen

Department of Electronic Engineering, China Institute of Technology & Commerce
No. 245, Sec. 3, Yen-Chiu-Yuan Road, Taipei, Taiwan

ABSTRACT

In this paper, a hybrid Proportional-Integral-Derivative (PID) controller that combining the Ziegler-Nichols PID controller with the grey prediction PID controller is proposed. The fuzzy gain scheduler is constructed to integrate these two controllers. Different characteristics of the employed controllers have been appropriately acquired. The great advantage of the proposed control architecture is that the parameters of PID controllers do not need to adapt. Furthermore, the design scheme provides an easy way to design the PID controller. The goal of the first Ziegler-Nichols PID controller is designed with fast response. Usually, it can be obtained after using the Ziegler-Nichols tuning algorithms. The second grey prediction PID controller is operated in slow response. It can be easily achieved through scheduling the system output. According to the employing different characteristics of controllers, the fuzzy gain scheduling has been successfully applied to emulating these two controllers to take care of the transient and steady state's performance simultaneously under the situations of unchanging the parameters of PID controller. Simulation results exhibit the superiority of the proposed method over the conventional ones.

Keywords: Fuzzy gain scheduling, Grey systems, Grey prediction, PID control

INTRODUCTION

Most of the control techniques implemented in industrial processes employ PID controller. There are two reasons why nowadays it is still the majority in industrial processes. The first reason is that its simple structure and the well-known Ziegler-Nichols tuning algorithms have been developed [1-2]. The second reason is that the controlled processes in industrial plant almost can be controlled through the PID controller [3-4]. However, the conventional PID controller design usually needs to retune the parameters (proportional gain, integral time constant and derivative time constant) mutually by a skilled operator. In particular, in order to improve its performance the fuzzy set theory is incorporated to tune the parameters of PID controller [5-9]. Generally speaking, the means of fuzzy tuning provide effective methods in the PID controller design. However, the established tuning rules are deeply based on someone who has rich knowledge/experiences. Usually, it is not an easy task to construct the rules for these parameters. Furthermore, the relations between the parameters are intertwined. Therefore it needs to take much time to characterize the parameters of PID controller. In this paper, an easy but effective control architecture of PID controller that integrating the well-known Ziegler-Nichols PID controller with grey prediction controller is introduced. In order to compensate the characteristic of original controller, the predicted system output feeds into the PID controller. Essentially, different system performance can be obtained if using the different prediction step. We can first take several PID controllers that have different performances in distinct operation conditions. Then, according to these local performances of PID controllers, the fuzzy gain scheduling technique is appended to determine the contribution/gain of every PID controller. Noted that the parameters of PID controller are unchanged after applying Ziegler-Nichols tuning process.

Grey prediction method initially presented by Prof. Deng [10-11]. The great advantage of grey prediction is that it only needs several data to develop a grey model. Instead of considerable data and well behavior of distribution are used by conventional prediction methods, only at least four data are needed in grey prediction. Thus it is possible to apply it to the requirement of real-time control systems. The grey prediction method employed in control system has some basic characteristics described as follows [12-14]. The large prediction step using in prediction PID controller usually causes worse transient response, but small prediction step yields fast response. For avoiding change the parameters of PID controller, the fuzzy gain sched-

uling is introduced. According to the scheduling of the system output via grey prediction, the fuzzy gain scheduler is constructed by some gain tuning rules to emulate the amount control signal of the selected controllers. Therefore, how to integrate the grey prediction PID controllers with the fuzzy gain scheduling is the core of designing the control system.

On the basis of linear control theory, the gain scheduling has been widely used in controlling systems whose dynamics change with the operating (operation) conditions [15-19]. Basically, it normally requires to know a conventional model of the nonlinear system and the partition of the state space under control. Then the well developed linear controllers are designed with respect to the linearized system at each operation conditions. In order to avoid the abruptness of the controller's parameters while across the transition region, the fuzzy gain scheduling is proposed to smooth the parameters of controllers. In this study, it is assumed that the linear model has been already achieved for convenience, and the controllers with distinct characteristic are well designed. Then, the aim of the proposed fuzzy gain scheduling is utilized to determine the weight of the each controller.

GREY PREDICTION

The basic concept of grey prediction employ finite data (at least 4 data) to construct the grey prediction model. The first step of grey model (GM) accumulates the selected data such that the accumulated data sequence is more regular than original data sequence. According to the solution of the grey differential equation, the prediction value for the regulated data can be obtained through a giving prediction step. As we obtain the prediction value, the inverse accumulated operation is applied to getting the prediction value of the original data. More explanation about them reader can refer to [20-21] for details. A brief introduction is given as follows.

Suppose the system output is interested to predict. The measured data in time sequence are denoted as

$$x_o^{(0)} = \{x_o^{(0)}(1), x_o^{(0)}(2), \dots, x_o^{(0)}(n)\}. \quad 1.$$

where n is the sample size of the record data. The original data is usually modified by mapping operation for establishing the efficiency of GM model.

$$x^{(0)}(k) = \exp(\gamma \cdot x_o^{(0)}(k)), \quad k = 1, 2, \dots, n \quad 2.$$

where 'exp' denotes exponential operation.

After the operation of data mapping, the accumulated generating operation (AGO) is defined as

$$x^{(1)}(k) = \sum_{i=1}^k x^{(0)}(i), \quad k = 1, 2, \dots, n. \quad 3.$$

or it can be rewritten in the following notation

$$x^{(1)}(k) = AGO[\exp(\gamma \cdot x_o^{(0)}(k))], \quad k = 1, 2, \dots, n \quad 4.$$

Obviously, the $x^{(1)}$ data sequence becomes more regular than the original data sequence, and it exhibits strictly increased data sequence. The objective of grey modeling (GM) is to find out the developing law of the regulated data sequence via the differential equation. First, define $z^{(1)}$ as the data sequence obtained by the following MEAN generating operation to $x^{(1)}$,

$$z^{(1)}(k) = \frac{1}{2} [x^{(1)}(k) + x^{(1)}(k-1)], \quad k = 1, 2, \dots, n. \quad 5.$$

Then the equation

$$x^{(0)}(k) + ax^{(1)}(k) = b. \quad 6.$$

is called a grey differential equation of GM and the whitening differential equation corresponding to the grey differential equation is

$$\frac{dx^{(1)}}{dt} + ax^{(1)} = b \quad 7.$$

where a and b are the developing coefficient and the grey input, respectively. In order to find out the solution of (6), the coefficients should be determined in advance. Define the grey parameter vector $g = [a \ b]^T$, then by the least square method, we have

$$g = \begin{bmatrix} a \\ b \end{bmatrix} = (A^T A)^{-1} A^T B \quad 8.$$

where

$$A = \begin{bmatrix} -0.5(x^{(1)}(1) + x^{(1)}(2)) & 1 \\ -0.5(x^{(1)}(2) + x^{(1)}(3)) & 1 \\ \vdots & \vdots \\ -0.5(x^{(1)}(k-1) + x^{(1)}(k)) & 1 \end{bmatrix}, \quad B = \begin{bmatrix} x^{(0)}(2) \\ x^{(0)}(3) \\ \vdots \\ x^{(0)}(k) \end{bmatrix}. \quad 9.$$

Therefore, based on the solution of (6), the GM model with respect to the data sequence $x^{(1)}$ is given by

$$\hat{x}^{(1)}(k + p_s) = (x^{(0)}(1) - \frac{b}{a})e^{-a(k+p_s-1)} + \frac{b}{a}, \quad k = 1, 2, \dots, n. \quad 10.$$

where p_s denoted as prediction step size or p -step ahead.

The inverse accumulated generating operation (IAGO) is used to estimate the value of (2) at p -step ahead

$$\hat{x}^{(0)}(k + p_s) = \hat{x}^{(1)}(k + p_s) - \hat{x}^{(1)}(k + p_s - 1). \quad 11.$$

Similarly, the inverse data mapping operation should be applied to $\hat{x}^{(0)}$ for obtaining the prediction value of the measurement data sequence at p -step ahead

$$x_o^{(0)}(k + p_s) = \frac{1}{\gamma} \ln(\hat{x}^{(0)}(k + p_s)). \quad 12.$$

In summary, the overall operation of the grey prediction can be simplified as

$$\hat{x}_o^{(0)}(k + p_s) = \frac{1}{\gamma} \ln[IAGO \cdot GM \cdot AGO(\exp(\gamma \cdot x_o^{(0)}(k)))] \quad 13.$$

GAIN SCHEDULING

Consider a nonlinear system governed by the following dynamic equation

$$\dot{\tilde{x}} = f(\tilde{x}(t), u(t)), \quad y = x_1 \quad 14.$$

where $\tilde{x} = [x_1, x_2, \dots, x_n] \in R^n$ is an $n \times 1$ state vector, u is the control input, and y is the system output. Assume that there exists a known family of operation point, say $(\tilde{x}_{op}^i, u_{op}^i)$, $i = 1, 2, \dots, k$. Then, the linearization around each operation point in the state space of error results in

$$\dot{\tilde{e}} = A_i \tilde{e} + b_i e_u, \quad 15.$$

where $\tilde{e} = \tilde{x} - \tilde{x}_{op}^i$, $e_u = u - u_{op}^i$, $A_i = \frac{\partial f}{\partial \tilde{x}} \Big|_{\tilde{x}_{op}^i, u_{op}^i}$, and $b_i = \frac{\partial f}{\partial u} \Big|_{\tilde{x}_{op}^i, u_{op}^i}$. The first step of gain scheduling technique is to design the controller with respect to each linearization model. Then each control law can be represented as

$$u_i = C_i(\tilde{e}), \quad 16.$$

where $C_i(\cdot)$ is the i th controller for the i th linearized model. Essentially, the control laws of controllers are changed with different operation points. The parameters regulation of controllers between the operation

points is the issue of gain scheduling. Based on the technique of gain scheduling, the fuzzy gain scheduling is adopted to tune the gain of PID controllers.

HYBRID PID FUZZY GAIN SCHEDULER DESIGN

The PID controller generates a control law $u(t)$ based on the closed-loop error, which can be formulated as

$$u(t) = k_p \left[e(t) + \frac{1}{T_i} \int e(t) dt + T_d \frac{de(t)}{dt} \right] \quad 17.$$

where k_p , T_i , and T_d are proportional gain, integral and derivative time constants, respectively, and $e(t)$ is the error between the reference $r(t)$ and the system output $y(t)$. The discrete time expression for PID controller is represented as

$$u(k) = k_p \left(e(k) + \frac{T_s}{T_i} \sum_{j=1}^k e(j) + \frac{T_d}{T_s} \Delta e(k) \right) \quad 18.$$

where $\Delta e(k)$ is the error rate between the error of the present time instant and the last time instant, $\Delta e(k) = e(k) - e(k-1)$, T_s is the sampling period for PID controller. The criteria to search the parameters of PID controller are based on evaluating the stability limits of the system. The process of tuning schemes can be described as follows: (1) The integral and derivative terms of PID controller are initially taken out of the system and remain the proportional term only, (2) The proportional gain k_p is increased until continuous oscillations are observed (marginally stable), (3) The corresponding gain k_p as k_x (ultimate gain) and the period T_x (ultimate period) of the oscillation are recorded, (4) Finally, the Ziegler-Nichols tuning of the PID controller are suggested as $k_p = 0.6k_x$, $T_i = 0.5T_x$, and $T_d = 0.125T_x$.

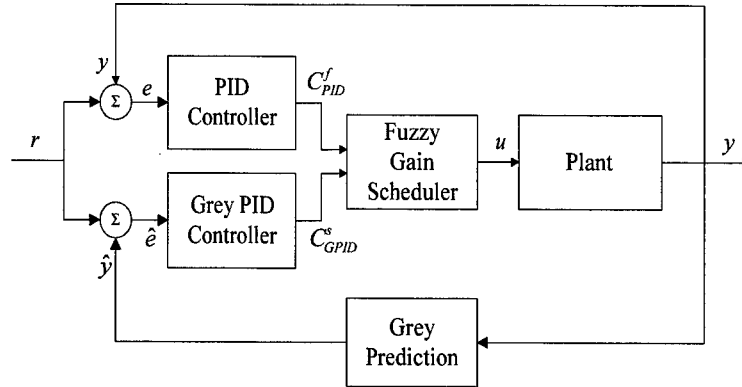


Fig. 1. The block diagram of the proposed PID controller

The block diagram of the proposed control scheme is depicted in Figure 1, where the grey prediction is used to estimate the system output so that the obtained error is always different from the error of present state. The controller is composed of the following three parts: a PID controller, a grey prediction PID controller and a fuzzy gain scheduler. The input variable of PID controller is $e(t)$, whereas, the grey prediction PID controller is $\hat{e}(t)$. The prediction error and the prediction error rate of system at p th-step ahead are respectively defined by

$$\hat{e}(k+p) = r - \hat{y}(k+p), \quad \Delta \hat{e}(k+p) = \hat{e}(k+p) - \hat{e}(k+p-1) \quad 19.$$

where r is the reference signal, \hat{y} is the system prediction output. Then, according to (18) the control law of grey prediction PID controller is

$$u(k) = k_p \left(\hat{e}(k+p) + \frac{T_s}{T_i} \sum_{j=1}^k \hat{e}(j+p) + \frac{T_d}{T_s} \Delta \hat{e}(k+p) \right) \quad 20.$$

Conventionally, the design of grey prediction controller is almost demonstrated to how tune the prediction step of the grey prediction [12-14]. From their studied, we can find the fact that when the large prediction

step used in grey prediction the poor rise time is happen, while using small prediction step, the response is similar but slightly better than that without using grey prediction, that is, it usually obtains fast response (smaller rise time), large overshoot and more settling time. Therefore, the fuzzy rule base [12], stochastic learning [13] and switching algorithms [14], are used to change the prediction step dynamically, such that the performance of control system can be better than that of original controller with appending fix grey prediction step and without grey prediction. In this study, we take the characteristics of grey prediction and according to the schemes of gain scheduling under the situation that fix prediction steps and the parameters of PID controller. A fuzzy gain scheduling is proposed to emulate each controller's contribution and to determine the gain of each controller, such that the performance can be considerably improved. It is well-known that the faster response, lower overshoot and shorter settling time is usually required in most control applications. However, these basic requirements are trade-off. We exploit two of these PID controllers. The first one is the pure Ziegler-Nichols PID controller and aimed at fast response. The second one is the grey prediction PID controller with smooth transient response. The selected controllers can be respectively achieved through the Ziegler-Nichols tuning process and large prediction step. The fuzzy gain scheduling is incorporated to determine the weights of selected controllers in order to achieve the control goals with the faster response, lower overshoot, and shorter settling time simultaneously.

To begin with, we make the assumption that there are two PID controllers already achieved and denoted as C_{PID}^f (fast Ziegler-Nichols PID controller) and C_{GPID}^s (slow grey prediction PID controller). Moreover, we assume that there are two same linearized systems for the reason of simplification. Thus gain scheduling can be involved in our control scheme. The gain rule base of fuzzy gain scheduling to the controllers are formed as

$$R_i : \text{IF } |e| \text{ is } A_i \text{ then } u \text{ is } \alpha_i C_{PID}^f + \beta_i C_{GPID}^s, \quad i=1,2,\dots,n. \quad 21.$$

where R_i is the i th rule. According to the defuzzification of weighting average, the control law is

$$u = \frac{\sum_{i=1}^n \mu_{A_i}(|e|) [\alpha_i C_{PID}^f + \beta_i C_{GPID}^s]}{\sum_{i=1}^n \mu_{A_i}(|e|)} = \sum_{i=1}^n \mu_{A_i}(|e|) [\alpha_i C_{PID}^f + \beta_i C_{GPID}^s] \quad 22.$$

where $|\cdot|$ is absolute value, and the membership functions of IF-part are defined in uniformly and symmetrically distributed over the universe of discourse so that $\sum_{i=1}^n \mu_{A_i}(|e|) = 1$. Equation (21) implied that the control law of hybrid PID controller can be changed by $\alpha_i, \beta_i, i=1, \dots, 2$, and the membership functions of IF-part. The determination of these parameters could be appropriately selected by the rule of thumb. The design procedure is summarized as follows. Step 1: Determine the parameters of Ziegler-Nichols PID controllers, Step 2: Append the grey prediction to Ziegler-Nichols PID controllers, Step 3: Choose the controllers, one with fast transient response, another with slow response, Step 4: Construct the rules of fuzzy gain scheduler, and determine the amount of control signal according to (22), Step 5: Perform the hybrid PID controller.

SIMULATION RESULTS

In this section, to exhibit the better performance of the proposed approach, two simulation examples are used to verify. One is third-order process, and another is fourth-order process. They can be respectively represented as the following transfer functions

$$G_1(s) = \frac{4.228}{(s+0.5)(s^2+1.64s+0.456)} \quad 23.$$

$$G_2(s) = \frac{27}{(s+1)(s+3)^3} \quad 24.$$

The step responses for the processes (23) and (24) are used in the following strategies for the reason of comparison. (a) apply conventional Ziegler-Nichols PID controller; (b) append the grey prediction to Ziegler-Nichols PID controller with the small prediction step and the large prediction step to characterize the effects of grey prediction PID controllers; (c) apply fuzzy gain scheduling to the PID controllers to determine the weight of the selected controllers. The sampling time interval for computer simulation is given as

0.01sec. The ultimate gain and the period of Ziegler-Nichols PID controller are determined as, $k_x = 3.6$ and $T_x = 2.1s$ for process (23), $k_x = 4.88$ and $T_x = 2.55s$ for process (24).

The rules and the membership functions of fuzzy gain scheduler are respectively given in Table 1 and Figure 2. As usual, the meanings of the fuzzy sets are, respectively, PL for Positive Large, PM for Positive Medium, PS for Positive Small and ZR for Zero.

Table 1. The rules of the fuzzy gain scheduler

$ e $	PL	PM	PS	ZR
u	$\alpha_1 C_{PID}^1 + \beta_1 C_{PID}^2$	$\alpha_2 C_{PID}^1 + \beta_2 C_{PID}^2$	$\alpha_3 C_{PID}^1 + \beta_3 C_{PID}^2$	$\alpha_4 C_{PID}^1 + \beta_4 C_{PID}^2$

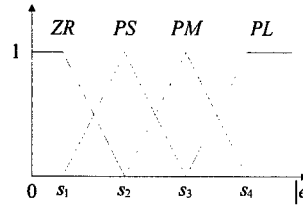


Fig. 2. The membership functions of the IF-part of the fuzzy gain scheduler

The step responses using conventional Ziegler-Nichols PID controller and grey prediction PID controllers with small prediction step (3-step) and large prediction step (150-step) are shown in Figure 3. As can be seen from Figure 3, the small prediction step used by grey prediction PID controller is almost the same as the Ziegler-Nichols PID controller. The more rise-time is required if the large prediction step is used in grey prediction PID controller. We use these two selected controller, one is the Ziegler-Nichols PID controller, another is the grey prediction PID controller to construct the fuzzy gain scheduler based on the fuzzy rule manner. The parameters of fuzzy gain scheduler for both processes are determined as $[s_1, s_2, s_3, s_4] = [0.25, 0.5, 0.75, 1]$, $[\alpha_1, \alpha_2, \alpha_3, \alpha_4] = [1, 0.95, 0.9, 0.9]$ and $[\beta_1, \beta_2, \beta_3, \beta_4] = [0.2, 0.5, 0.8, 1]$. The step responses of the hybrid PID controller with fuzzy gain scheduling for processes (23) and (24) are respectively shown in Figure 3 (a) and (b). Obviously, the proposed control approach has better results than that of the conventional one.

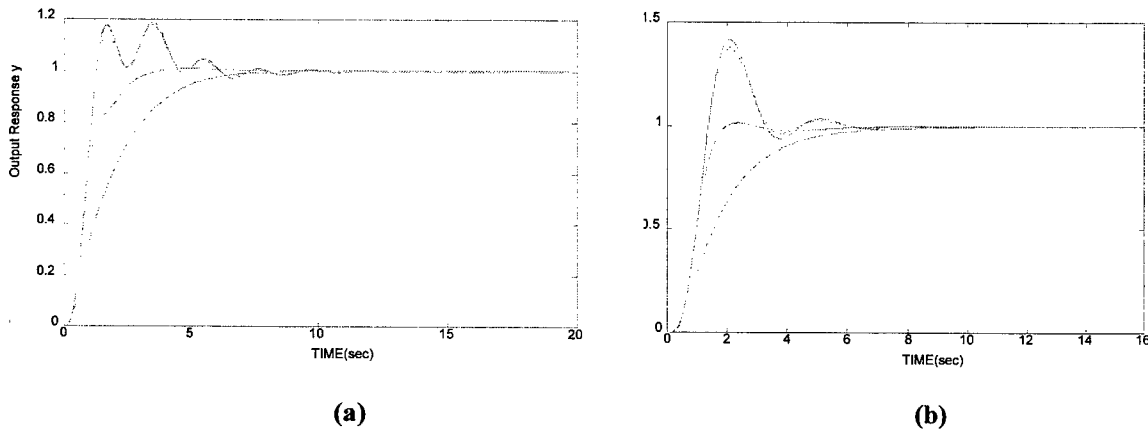


Fig. 3(a). The step responses of the third-order process with Ziegler-Nichols PID controller (dotted line), grey prediction PID controller at 3-step ahead (dash dotted line) and at 150-step ahead (dashed line) and the proposed hybrid PID controller with fuzzy gain scheduler (solid line). **(b)** The step responses of the fourth-order process with Ziegler-Nichols PID controller (dotted line), grey prediction PID controller at 3-step ahead (dash dotted line) and at 150-step ahead (dashed line) and the proposed hybrid PID controller with fuzzy gain scheduler (solid line).

CONCLUSIONS

A hybrid PID controller has been well designed through fuzzy gain scheduling. Instead of tuning the parameters of PID controller used by conventional approaches, the technique of fuzzy gain scheduling is employed to determine the gain of well-designed grey prediction PID controllers. The great advantage of the proposed approach is that the parameters of original Ziegler-Nichols PID are unchanged through system operation. Both of the transient and steady state response are considered simultaneously, which is not easily achieved if the conventional approaches are applied. Considering the characteristic of grey prediction and the essential concepts of gain scheduling, a set of heuristic rules of fuzzy gain scheduling has been constructed to determine the amount control signal depending on error. Basically, the proposed approach provides an effective way to construct the PID controller. From the simulation studies, we can find that the hybrid PID controller is superior to the conventional Ziegler-Nichols PID controller.

ACKNOWLEDGMENT

This research was supported by the National Science Council, Republic of China, under contract NSC 87-2218-E-157-004.

REFERENCES

1. G.F. Franklin, J.D. Powell, E.N. Abbns, 1988. Feedback control of dynamic systems, Addison Wesley.
2. B.C. Kuo, 1987. Automatic control systems, 5th ed. Englewood Cliffs, NJ:Prentice-Hall.
3. C.C. Hang, K.J. Astrom, W.K. Ho, 1991. Refinements of the Ziegler-Nichols tuning formula, Proc. IEE Pt. D, 138, 111-118.
4. Z.Y. Zhao, M. Tomizuka, S. Isaka, 1993. Fuzzy gain scheduling of PID controllers, IEEE Trans. on Syst., Man, and Cybern., 23(5), 1392-1398.
5. W. Pedrycz, J.F. Peters, 1997. Hierarchical fuzzy controllers: fuzzy gain scheduling, IEEE Inter. Conf. on SMC, 1139-1143.
6. S.Z. He, S. Tan, F.L. Xu, P.Z. Wang, 1993. Fuzzy self-tuning of PID controllers, Fuzzy Sets and Systems, 56, 37-46.
7. Q.P. Ha, M. Negnevitsky, F. Palis, 1997. Cascade PI-controllers with fuzzy tuning, IEEE Inter. Conf. on Fuzzy Systems, 361-366.
8. S. Tzafestas, N.P. Papanikolopoulos, 1990. Incremental fuzzy expert PID control, IEEE Trans. on Industrial Electronics, 37(5), 365-371.
9. V. Pauli, N.K. Heikki, 1995. Fuzzy logic in PID gain scheduling, 3rd European Congress on Intelligent Techniques and Soft Computing, 2, 927-931.
10. J.L. Deng, 1982. Control problems of grey system, System and Control Letters, 1(5), 288-294.
11. J.L. Deng, 1989. Introduction to grey system theory, J. of Grey System, 1(1), 1-24.
12. C.M. Hong, S.C. Lin, C.T. Chiang, 1995. Control of dynamic system by fuzzy-based grey prediction controller, J. of Grey System, 7(1), 23-44.
13. J.Y. Chen, 1996. Tuning of prediction step in grey prediction controllers using stochastic learning, J. of Grey System, 8(4), 337-357.
14. C.C. Wong, C.C. Chen, 1997. Switching grey prediction PID controller design, J. of Grey System, 9(4), 335-350.
15. W.J. Rugh, 1990. Analytical framework for gain scheduling, IEEE control Syst. Mag., 11(1), 74-84.
16. D.A. Lawrence, W.J. Rugh, 1995. Gain scheduling dynamic linear controllers for a nonlinear plant, Automatica, 31(3), 381-390.
17. R. Palm, U. Rehfuess, 1997. Fuzzy controllers as gain scheduling approximators, Fuzzy Sets and Systems, 85, 233-246.
18. S. Tan, C.C. Hang, J.S. Chai, 1997. Gain scheduling: From conventional to neuro-fuzzy, Automatica, 33(3), 411-419.
19. L. Cheng, E.F. Thomas, 1997. Real time control of a water-gas shift reactor by a model-based fuzzy gain scheduling technique, J. of Process Control, 7(4), 239-253.
20. J.L. Deng, 1992. The essential methods of grey systems, HUST Press, Wuhan, in Chinese.
21. H.C. Lu, 1996. Universal GM(1,1) model based on data mapping concept, J. of Grey System, 8(4), 307-319.

Holonically Object-Oriented System

Shigeki Sugiyama

Gifu Industry and Technology Research Center,
47 Kataoyobi, Kasamatsu-Cho, Hashima-Gun, Gifu-Ken, Japan.
Email: sugiyama@vsl.gifu-u.ac.jp, sugiyama@gifu-irtc.go.jp

ABSTRACT

In this paper, we introduce the concept of Holonically Object Oriented Systems. Nowadays there are many more complicated things than ever before in the world, waiting to be controlled intelligently in order to improve the production rate in a factory; make things more clear in a complex system; or get help in terms of analyzing a system, etc. To accomplish this, we have tried to increase our understanding, accuracy and precision of the target system to be controlled. After obtaining the required information, we are ready to control a system for many different purposes. But often this approach can complicate a problem further, which then becomes more time-consuming because of an increase in system size, resulting in comparatively low robustness. This can be caused by: a lack of a flexibility against sudden changes in the behavior of a system; by giving too much redundant attention to a particular aspect of the system and; by a lack of intelligence.

PRESENT SITUATION

In the past, a production system was considered as an integration of workers. Then the definition became an integration of machines; then we developed semi-automatic machines which provided lower labor costs and more efficient equipment; then came intelligent production systems; and so on. The world has been getting smaller and smaller because of the development of information technologies and high-speed transportation systems. So, in the case of production of goods, it was natural that the production base has shifted to places which can offer good and reasonable labor costs. But naturally, this phenomenon soon ends up with rising labor costs as these new places grow, change and mature. So now, we must turn to the age of agile production with respect to cost and quality with more complex, large and rapidly changing systems than ever before. In addition to this situation, a company cannot survive without global-expansion in relation to production, workers, information and information technology. If these factors are integrated, it is possible to produce a competitive product.

In such situations, it is difficult to maintain system stability with respect to disturbances, to provide high adaptability and flexibility to change, and to sustain effective production. So, in this atmosphere of requiring more and more effective production, we must have a mechanism in the organization that will react intelligently to:

- unpredictable modification of products,
- increasing numbers of product variations,
- frequent model changes,
- shorter product life-time,
- constant production line.

On the other hand, we have problems:

- system becomes rigid due to expansion in size,
- limits of automation are reached, diminishing returns,
- wastage of human skills and talent,
- adaptability to change in the environment is low,
- inability to guarantee future progress.

So from the above, it is obvious and important to have flexible behaviour to express the system correctly and accurately. Secondly, we need to reduce redundant behaviour and provide effective production.

To address these problems, we introduce the Holonically Object Oriented System. To deal with the above problems, it is necessary to generalize the SYSTEM to be considered. So, in the next section, a general expression of the SYSTEM is considered.

GENERAL CONCEPT OF THE SYSTEM

Systems in General

Everything that has some sort of "Organism" or "Pattern of Behavior" or "Regulated Phenomenon" or "Relation" inside or within itself, can be defined as a "SYSTEM" in which the behaviour or interaction phenomena of creatures with their environment can be well-defined by equilibrium theory, the laws of dynamics and/or the laws of thermodynamics [1], and the behaviour or interaction phenomena of things amongst themselves is of utmost interest [2].

According to "The Structure and Function of Organization" by J. Feibleman and J.W. Friend [3], the study of organizations must be approached from two viewpoints - in either a static or dynamic way. Static treatment views an organization as being independent of its environment and therefore isolated from interaction problems with other organizations. Dynamic treatment views an organization as being dependent, to some extent, upon its environment and therefore, it interacts with other organizations. As mentioned above, our interest is focused on the behavior of an organization itself, so static treatment is our only concern. In one sense, the system we treat is an open system which attains a time-independent state wherein the system remains constant as a whole and constant within each of its phases or parts, although there is a continuous flow of component materials [4]

The General Concept of a Static Organization

The Basis of Organization

J. Feibleman and J.W. Friend define the Static approach as; "in treating structure, we first examine the organization itself as a whole. The whole obviously analyses into parts. These parts themselves have parts, which we shall term, subparts.

Thus there are two levels of analysis:

1. Wholes(from which an analysis is made); and
2. Parts and Subparts.

So from this, we can have the following conjecture.

Conjecture 1:

Any static Organization can be recognized as a Whole and the Whole can be treated as a static System. The System has parts, and the parts consist of subparts. The subparts have their sub-subparts, and so on.

When we think about this world, we can easily see the following phenomenon:

At one time, the whole is the main function of a System, but at another time, each part or subpart may function as the main feature of the system.

So this leads to the next conjecture:

Conjecture 2.

In a system, at one moment, "HOLOS" can be the main function that explains the whole behavior of the system, and at another moment, each "ON" can be the main function which explains the whole behavior of the system. So the phenomenon of "HOLOS-ON ? HOLON" is a function which can be found in every system.

Now, lets describes the elements of relationships [4]. There is an important factor in the analysis of organizations that we may temporarily describe as the ways in which parts exist in combination with other parts to form the structure of the whole. We refer here to the kinds of relationships between parts and the ways in which parts combine. The elements of relationships which exist between parts of an organization form a group of relationships as follows:

1. Transitivity: If we relate two parts to a middle one, this relates the extreme parts to each other.
2. Connectivity: This is the relationship of two parts without mediation of a third part.
3. Symmetry: This is a relationship between two parts which is the same in both directions.
4. Seriality: This is a relationship that is transitive, asymmetrical, and connected.
5. Correlation: (one-many, one-one, many-one, many-many) is a relationship between two series such that for every part of one series there is a corresponding part in the other series and no part in either series is without a corresponding part in the other.
6. Addition: This relationship joins parts together so as to increase their number within a part..
7. Multiplication: This relationship joins parts so as to involve them with each other.
8. Association: This is a relationship which is commutative and connected.
9. Distribution: This is a relationship which is commutative and intransitive.
10. Dependence : This is a relationship in which existence of one part is conditioned by some other part.

Rules of Organization

We now have the basis of organization, i.e., the whole, parts and subparts, and we also have elements of relationships between parts. But these are insufficient to define or determine any given organization. So in addition, we need certain rules in which, parts and their relationships are constitutive of organizations.

1. Structure is the sharing of subparts between parts.
The linkage of parts is accomplished by means of common subparts and not by mere juxtaposition or external linkage. The joining of two parts is effected by a subpart which they hold in common, and this is the basis upon which all structures are constituted.
2. Organization is the one controlling order of structure.
It is not the facet of linkages but rather, the principle under which all linkages fall into one controlling order, which makes an organization.
3. One additional level is needed to constitute an organization besides its parts and subparts.
No specific number of subparts and parts constitutes an organization, which essentially, is a property of the whole. The organization is one level above its analytic parts and subparts, and thus the whole must involve an additional level.
4. In every organization there must be a serial relation.
The serial relation is essential in every organization. Other relationships may and usually do exist but they are unnecessary to constitute an organization. In analyzing every whole, there must be a controlling relationship which is asymmetrical, transitive, and connective.
5. All parts are shared parts.
There is nothing in an organization except parts which have subparts in common. An item which is not shared by parts is extraneous to the organization and not a part.
6. Things in an organization related to parts of the organization are themselves parts of the organization.
Anything in an organization which is related to part of that organization, by virtue of that relationship, is itself part of the organization, and not a foreign body.
7. Things in an organization related to related parts are themselves parts of the organization.
Sometimes there are things in an organization which are not related to any single part of the organization, but which are related to two or more parts taken together.
8. The number of parts and number of their relationships constitutes complexity.
The number of parts and their relationships, i.e., subparts, constitute the complexity of an organization. This role and the yardstick of integrality, or kinds of static organization form a pair of criteria. Complexity is reduced to a mere matter of counting parts and subparts.

Kinds of Organization

The kinds of organization constitute degrees of integrality as follows:

1. Agglutinative
The governing relation is aseriality, where parts have intransitivity, connectivity and symmetry.
2. Participative
The governing relation is seriality. Participative organizations subdivide into three kinds;
 - a. Adjunctive
The governing relationship is symmetrical independence. The sharing of subparts is not necessary to either of the parts. Parts can survive their separation.

b. Subjective

The governing relationship is asymmetrical dependence. The sharing of parts is necessary to one of the parts but not to both.

c. Complementary

The governing relationship is symmetrical dependence but in this case, the sharing of parts is necessary to both of the parts. Neither part can survive separation.

In the next section, by using the above terminology, the SYSTEM will be defined.

THE CONCEPT OF SYSTEM

Generally speaking, we can define the System of any Organization with conjecture 1 above as follows:

An organization has a whole and its parts. Parts have their subparts. And this phenomenon goes on and on into inside-subparts. So, we can define a system as follows:

Definition 1.

$$\begin{array}{ll}
 U = \{U|U\} & U : \text{Universe} \\
 P = \{P, p|p? ? P, ? P? U\} & P : \text{Parts} \\
 SP_0 = \{SP_0, sp_0 |? sp_0 ? ? SP_0, ? SP_0 ? P\} & p : \text{Elements} \\
 SP_1 = \{SP_1, sp_1 |? sp_1 ? ? SP_1, ? SP_1 ? SP\} & SP : \text{Subparts} \\
 \dots & sp : \text{Elements in Subparts} \\
 SP_n = \{SP_n, sp_n | sp_n ? ? SP_n, ? SP_n ? SP_{n-1}\} &
 \end{array}$$

These parts and subparts, etc. have a relationship among themselves as described in the Basis of Organization section above. These are "transitivity, connectivity, symmetry, seriality, correlation, addition, multiplication, commutation, association, distribution and dependence". These relationships are categorized into Relations, Functions, and Connections between parts and subparts as follows:

RELATION	transitivity, connectivity, symmetry, seriality
FUNCTION	addition, multiplication, commutation, association, dependence
CONNECTION	correlation

Generally speaking, it can be said that an organization will contain factors of RELATION, FUNCTION, and CONNECTION. Furthermore, these factors perform the following tasks:

RELATION [R] specifies how to transfer,
 FUNCTION [F] specifies a quantity to be transferred,
 CONNECTION [C] specifies a route to transfer.

So from the facts of a System (Definition 1.) and an organization as mentioned above, we can define a System of Organization (SO) as:

Definition 2.

$$SO(X) = \{X | [R], [F], [C], U, P, SP_1, \dots, SP_n\}$$

$$SO(U, P, SP, \dots) = \{[R], [F], [C]\}$$

where $[R] = f(t, c, sy, se)$ $t = \text{transitivity}, c = \text{connectivity}, sy = \text{symmetry}, se = \text{seriality}$
 $[F] = f(\text{add}, m, c, \text{ass}, d)$ $\text{add} = \text{addition}, m = \text{multiplication}, c = \text{commutation}, \text{ass} = \text{association}, d = \text{dependence}$
 $[C] = f(r)$ $r = \text{route}$

From the above, the elements of RELATION are specified by FUNCTION [F] and this information goes to with the route defined by CONNECTION [C]. So this can be rewritten as:

$$\begin{array}{l}
 U = f([R([F])])_{[C]} \\
 P = f([R([F])])_{[C]} \\
 SP_1 = f([R([F])])_{[C]} \\
 \dots \\
 SP_n = f([R([F])])_{[C]}
 \end{array}$$

So we can define a System of Organization as follows:

$$SO(X,Y,Z,...) = f(U(X), P(Y), SP_1(Z_1),...,SP_n(Z_n)) \quad 1.$$

THE CONCEPT OF A HOLONICALLY OBJECT ORIENTED SYSTEM

So far, we have defined a SYSTEM OF ORGANIZATION (SO), so now it is necessary to introduce a system which will solve the specific problem case.

HOLON at Present

Twenty five years ago, Arthur Keostler proposed "HOLON" to describe a basic unit of organization in biological and social systems. HOLON is a combination of the Greek word "holos", meaning whole, and the suffix "on" meaning particle or part. Keostler observed that in living organisms and in social organizations entirely self supporting, non-interacting entities did not exist. Every identifiable unit of organization, such as a single cell in an animal or a family unit in a society, compresses more basic units while at the same time forming a part of a larger unit of organization. A holon, as Keostler devised the term, is an identifiable part of a system that has a unique identity, yet is made up of sub-ordinate parts and in turn is part of a larger whole.

The strength of a holonic organization, (holarchy) is to enable construction of very complex systems that require improvment in efficient use of resources, highly resilient to disturbances, and adaptable to environmental change. All such characteristics are seen in biological and social systems [5].

Definition of HOLON

A holon is defined by the Consortium on Intelligent Manufacturing Systems as follows:

An autonomous and co-operative building block of a manufacturing system to transform, store and/or validate information and physical objects. The holon consists of an information processing part and often, a physical processing part. A holon can be part of another holon.

Autonomy is the capability of an entity to create and control execution of its own plans and/or strategies. Co-operation refers to the process by which a set of entities develops mutually acceptable plans and executes these plans. A Holarchy is a system of holons that can co-operate to achieve a goal or objective. The holarchy defines the basic rules for co-operation of the holons and thereby limits their autonomy.

Holonic manufacturing system (HMS) is a holarchy that integrates an entire range of activities from order-booking through design, production, and marketing to realize an agile manufacturing enterprise. Holonic attributes are features of an entity that make it a holon. The minimum set of attributes is autonomy and cooperation.

Definition of a HOLONIC SYSTEM of Organization

In the field of control we have been trying to expand robustness since a manufacturing system to be controlled has become so large and too complicated. So this causes problems of instability, of sudden stopping of the whole system, of loss of control, of reduced reliability of the system, and so forth. To reduce this problem, a lot of things must be introduced such as enhanced intelligence, etc. But it is often difficult to achieve a satisfactory result because modern control still has the idea that the we must have the ability to control the whole system all the time. When we think about human or living creatures, we only use necessary functions as they are required. Other functions are left vague or redundant until needed. We can also apply this idea to the control of our present situation to reduce complexity and increase system reliability. This is the, so called, HOLONIC approach.

In the section on System of Organization, we defined our system mathematically, so now let us try to construct the System of Organization with a Holonic approach.

Holonic System of Organization

HOLON was first used in the book called "The Ghost in the Machine", London, 1967 by Author Koestler in which he introduced the idea of "Self-Regulating Open Hierarchic Order (SOHO)". This can be expressed as "YANUS", as well [5]. SOHO has the following characteristics:

- Living creatures do not consist simply of sets of parts and are not uni-motivated chains.
- The whole of each creature branches into sub-wholes one after another autonomously and each have a hierarchy of multi-levels.
- Self-regulating Open Hierarchic Order or YANUS
- Autonomy and integrality
- Hierarchy and networking
- Regulation and targeting
- The hierarchy has its own ordering regulations

In other words, we can conclude that a HOLON has characteristics of: parts and subparts, fluctuations, self-organizing; entrainment; YANUS; rhythm; distributed. These characteristics are expressed as:

- Parts and subparts are defined mathematically as $\{? P? U\}$
- Fluctuation is defined by a probability or a distribution function.
- Self-organizing is U or P or SP or SP_0 or ... or SP_n .
- Entrainment is the behavior of mimicking the behavior of neighbors.
- YANUS is a behavior of $\{? P? U\}$.
- Rhythm is defined by an input rate or a dispatching rate or a production rate.
- Distributive is defined by:

$$SO(X) = \{X, [R], [F], [C], U, P, SP_1, \dots, SP_n\}$$

$$SO(U, P, SP, \dots) = \{[R], [F], [C]\}$$

Now if we think about a company and its problems, the above characteristics can give us an easy way to express a company logically.

Concept of HOLONIC SYSTEM

Here the properties identified above are defined mathematically with a few necessary modifications.

Definition 3.

- Parts and subparts are defined mathematically as $\{? P? U\}$

2.

This was stated in Definition 1.

Definition 4.

Fluctuation is defined by a probability or a distribution.

The human heart does not always beat at exactly the same pace but fluctuates dependent on many factors: stress, lack of oxygen, heat, cold, etc.. We find the same behavior in many organs of our body. Furthermore, fluctuations are observed in every creature in the world. So $[F]$ is defined by:

$$[F] = \{\text{regulated } [r] \text{ and fluctuating behaviors } [f], \text{ probabilistic } [p] \text{ and distributive actions } [d]\}$$

Definition 5.

Self-organizing is U or P or SP or SP_0 or ... or SP_n .

Every part has its own self-organizing properties and we can find this property in every organism and creature. This can be defined as:

$$SO(X) = \{X, [R], [F], [C], U, P, SP_1, \dots, SP_n\}$$

3.

Definition 6.

Entrainment is the behavior of mimicking behavior of neighbors. We behave this way autonomously often without self-recognition, and again, this behavior can be found in every animal. Entrainment can be defined as:

$$f([F])_m = \{f([F]_n) ? f([F]_{n-1}) ? f([F]_{n-2}) ? \dots ? f([F]_1)\}$$

4.

Definition 7.

YANUS is a behavior of $\{? P? U\} ? \{U ? ?\}$ as it is. Mathematically, this is the same as $U=P$, but in this case it has another meaning, that is to say, every function in a system has a moment in which it is the whole $[U]$ of the system and so, at some instance, the whole may become a part $[P]$.

Definition 8.

Rhythm is defined by an input rate (IR) or a dispatching rate (DR) or a production rate (PR).

$$[R] = \{\text{input rate, dispatching rate, production rate}\}$$

Definition 9.

Distributive is defined by:

$$SO(X) = \{X | X? [R], X? [F], X? [C], X? U, X? P, X? SP_1, \dots, X? SP_n\} \quad 5.$$

and

$$SO(U(X), P(X), SP(X), \dots) = \{X? [R], X? [F], X? [C]\} \quad 6.$$

Definition 10.

A HOLONIC SYSTEM is a system which has properties defined by Definitions 3 - 9.

Definition 11.

HOLONIC CONTROL is a method which treats the HOLONIC SYSTEM.

Here we restate the System of Organization by using the definitions above. By using the definitions 1, 2, 7 and 9 described above, we have a general expression of a System of Organization as shown below.

$$SO(X, Y, Z) = f[U(X), P(Y), SP_1(Z_1), \dots, SP_n(Z_n)] \quad 7.$$

$$U = f([R([F]))]_{[C]})$$

$$P = f([R([F]))]_{[C]})$$

$$SP_1 = f([R([F]))]_{[C]})$$

$$\dots$$

$$SP_n = f([R([F]))]_{[C]})$$

By using and combining definitions 1, 2, and 3, we can restate the system which has both required properties of a HOLONIC expressed as:

$$SO(X, Y, Z)_H = f[U(X), P(Y), SP_1(Z_1), \dots, SP_n(Z_n)]_H \quad 8.$$

This can be rewritten as

$$SO(X? Y? Z)_H = \{U_H = f([R([F]))]_{[C]})_H ? P_H = f([R([F]))]_{[C]})_H ? SP_1_H = f([R([F]))]_{[C]})_H ? \dots ? SP_n_H = f([R([F]))]_{[C]})_H\} \quad 9.$$

By using definition 4, input quantities or properties have a certain amount of tolerance. They can be defined by [r], [f], [p], and [d]. These attributes are restated as [F] as used in the above equation. By using definition 6, we can restate the system as;

$$SO(X, Y, Z)_m = f[U(X)_m, P(Y)_m, SP_1(Z_1)_m, \dots, SP_n(Z_n)_m] \quad 10.$$

$$U_m = f([R([F]_m)])_{[C]}$$

$$P_m = f([R([F]_m)])_{[C]}$$

$$SP_{m1} = f([R([F]_m)])_{[C]}$$

$$\dots$$

$$SP_{mn} = f([R([F]_m)])_{[C]}$$

So with the above results, we can have a theory of the HOLONIC SYSTEM.

Theory 1.

A HOLONIC SYSTEM is a system which has attributes of System Organization and can be defined by the equation shown below:

$$SO() = \{SO(X, Y, Z)_m ? SO(X? Y? Z)_H\} \quad 11.$$

[Proof]

This is obvious from definitions 1~9.

We can state the normal System of Organization as:

$$SO(X,Y,Z) = f[U(X),P(Y),SP_1(Z_1),...,SP_n(Z_n)]$$

$$U = f([R([F]))][C]$$

$$P = f([R([F]))][C]$$

$$SP_1 = f([R([F]))][C]$$

...

$$SP_n = f([R([F]))][C]$$

And we can restate the above by adding the HOLONIC System as:

$$SO(X,Y,Z)_m = f[U(X)_m, P(Y)_m, SP_1(Z_1)_m, ..., SP_n(Z_n)_m] \quad 12.$$

$$U_m = f([R([F]_m))][C]$$

$$P_m = f([R([F]_m))][C]$$

$$SP_{m1} = f([R([F]_m))][C]$$

...

$$SP_{mn} = f([R([F]_m))][C]$$

And so, we can state the HOLONIC System as:

$$SO(X? Y? Z)_H = \{U_H = f([R([F]))][C] \quad H? \quad P_H = f([R([F]))][C] \quad H? \quad SP_1 H = f([R([F]))][C] \quad H? \\ SP_n = f([R([F]))][C] \} \quad 13.$$

From these two equations we can easily get the following result simply by adding two attributes:

$$SO() = \{SO(X,Y,Z)_m? \quad SO(X? Y? Z)_H\} \quad 14.$$

[End Proof]

Holonically Object Oriented System

In the above, we have cleared up the idea of HOLON (HOLONIC SYSTEM and HOLONIC CONTROL) and its mathematical definition. The idea of HOLON can be applied to any organization or company. A company has many departments, and each department consists of sections, and a section consists of men or equipment or facilities. Each has connections to one another. Each section or department has its own autonomous mechanism and can be independent itself.

So a system of a company will be able to define as an exactly same entity as shown below;

$$SO(X,Y,Z)_m = f[U(X)_m, P(Y)_m, SP_1(Z_1)_m, ..., SP_n(Z_n)_m]$$

U, P, SP₁,...and SP_n are independent of one another autonomously and also, they have close relationships with one another. And when all goes well, they behave just like one unit. We do not care about detail inside. But when the system is in a trouble, a section or department which has the trouble comes into sight. Sometimes this section or department will behave as if it is the whole. This is exactly the same as the HOLON defined above. So again, it can be expressed by theory 1 shown above.

$$SO() = \{SO(X,Y,Z)_m? \quad SO(X? Y? Z)_H\}$$

What is more, we can say that the each element of SO() can be expressed as an Object since it has autonomous behaviour and it behaves as one system.

And we can let it have the functions stated below:

Unless there is any trouble in X or Y or Z, since these behave independently, there is no need to be conscious about each of them. Once trouble occurs, it is very easy to know where it happens. The function SO will be controlled by data from the part in trouble.

In the simplest case, we can define the transformed outputs of X_o, Y_o, and Z_o for 3 levels as:

1. as expected; averaged ? C
2. linearly goes up or down ? aX+b.

And we can define the situation of C as follows:

$$SOo() = \{SOo(Xo)H\} \quad 15.$$

In this situation it is necessary to be conscious about the whole itself.

The situation of $aX+b$ can be defined as one of:

$$SOo() = \{SOo(Xo, Yo, Zo)m\} \quad \text{or} \quad 16.$$

$$SOo() = \{SOo(Yo)H\} \quad \text{or} \quad 17.$$

$$SOo() = \{SOo(Zo)H\} \quad 18.$$

In this situation it is necessary to be conscious about the Parts in trouble.

So the HOLONIC SYSTEM described above is easily transformed into an Object Oriented System, which is called a Holonically Object Oriented System. The above can be defined as shown below;

Definition 12.

A Holonically Object Oriented System is defined by the expression below:

$$SOo() = \{SOo(Xo, Yo, Zo)m? SOo(Xo? Yo? Zo)H\} \quad 19.$$

where $Xo = X$ expressed as an Object Oriented System
 $Yo = Y$ expressed as an Object Oriented System
 $Zo = Z$ expressed as an Object Oriented System

Applications in General

From the above results, we can apply company control. A company is expressed by the holarchy below:

COMPANY					
Management		Factory		Transfer	Sales
Research	Managing	Plan	Production	simulation	sale
market strategy	Order Accept Data analysis	plan simulation	control (RTM)	transfer	information foresee

COMPANY consists of Management, Factory, Transfer, and Sale. Management consists of Research and Managing. Factory consists of Plan and Production. Transfer consists of simulation and transfer. Sale consists of sale, information, and foresee. Research consists of market and strategy. Managing consists of order, accept, and data analysis. Plan consists of plan and simulation. Production consists of control.

COMPANY = {Management, Factory, Transfer, Sale}
 Management = {Research, Managing}
 Factory = {Plan, Production}
 Transfer = {simulation, transfer}
 Sales = {sale, information, foresee}
 Research = {market, strategy}
 Managing = {order, accept, data analysis}

As described above, this is able to define a Holonically Object Oriented System as below:

$$SOo = \{COMPANY\}$$

In this situation every element of COMPANY can be expressed as an Object and each Object can have a close relationship one to another. All transformations as expressed in equations 15., 16., 17., and 18. can be used in any particular situation.

CONCLUSION

We have introduced the Holonically Object Oriented System. The following results were obtained.

1. A HOLON is a very suitable idea to express a huge, complex system.
2. Any system has a holarchy and can be constructed with HOLON.
3. Holarchy has a good nature to be treated with Object.
4. Holonically Object Oriented System has the nature of flexibility, adaptability to change, and transformability into any form which itself remains to exist as an independent entity.

REFERENCES

1. A. Angyal, 1941. Foundation for a Science of Personality, Havard University Press.
2. W. Koehler, 1938. The Place of Values in the World of Fact, Liveright.
3. J. Feibleman, J.W. Friend, 1945. The structure and function of organization, Philosophical Review, Vol.54, 1945.
4. L.von Bertalanffy, 1950. The theory of open systems in physics and biology, Science 111.
5. P. Valckenaers, H.Van Brussel, 1998. "Holonically Manufacturing Systems. Techn. Overview HMS consortium report
6. Taro Nawa, 1988. Holonic Management Revolution, Sohko-sha Publishing Co.
7. R. Saucedo, E.E. Schiring, 1968. Introduction to continuous and digital control systems, Macmillan Publishing Co.
8. A.G.J. MacFarlane, 1970. Dynamical system models, George G. Harrap & Co.,

Intelligent Instrumentation and Measurement

Preprocessing of Industrial Process Data with Outlier Detection and Correction

J.Tenner*, D.A.Linkens* and T.J.Bailey**

* Department of Automatic Control & Systems Engineering,
University of Sheffield, U.K.

Email: J.Tenner@Sheffield.ac.uk

** British Steel Engineering Steels U.K. Ltd

ABSTRACT

When constructing predictive models from process data using techniques such as Neural Networks, the validity of the data is very important. This paper presents some current methods of 'cleaning' data and proposes a structured method applied to a batch heat treatment application in the steel industry. The methodology highlights the use of expert knowledge throughout a project's evolution. The application of this data cleaning methodology to the heat treatment process is described, and a quantitative comparison is made of the performance of a neural network model by comparing the accuracy of its predictions before and after the correction of outlying points.

INTRODUCTION

It is now common for industrial processes to be associated with large amounts of data. Increasingly this data is being used to construct empirical models to describe the underlying process [1,2,3]. One problem with this technique is that the data may contain points that are incorrect for one reason or another. These are known as outlying data points or outliers.

A popular technique for devising such models is that of Artificial Neural Networks, whose supervised learning types include Multi-Layer Perceptrons (MLPs) and Radial-Basis Functions (RBFs). Neural networks are trained by repeatedly examining a large number of examples for a particular problem. There are two main features that affect the success of this technique when applied to an industrial process. The first is the purity of the data upon which the network is trained, i.e., its ability to represent the process truthfully, and the second is the generalisation of the trained network when applied to previously unseen data.

Work has been undertaken to model heat treatment data from a range of steel processing sites of British Steel Engineering Steels U.K. This has been done largely using MLP networks with varying architectures.

BACKGROUND

The Steel Process

The heat treatment process for which the predictive model has been constructed is concerned with a wide range of low- to mid-carbon engineering steels with a variety of alloying elements. Heat treatment consists of three main stages: austenitising, quenching and tempering. The final mechanical properties of the steel are influenced by the temperature and quench used during the hardening and tempering stages, together with alloying elements which are added to the steel when it is being cast. The interaction effects between these parameters are often non-linear. Parametric knowledge is insufficient to develop suitable models, thus the aim of the model is to predict mechanical properties developed by the process for a wide range of engineering steels from past process data.

Neural Networks

In 1943, McCulloch and Pitts introduced the idea of an artificial neuron to process data [4]. In the 1950s, this work was advanced by arranging neurons in layers. Although learning rules to cope with multiple layers of perceptrons were not developed until later, this work formed the basis of MLPs used today.

A basic MLP is shown in Figure 1. It consists of a series of layers of neurons, namely the input, hidden and output layers. Each neuron in one layer is fully-connected to all neurons in the next layer. Each connection has a strength or weight associated with it. A network learns or is trained by modifying these weights. The result is a mapping from the input to the output layer via the hidden layers. Each neuron also has a bias value associated with it so as to provide an additional degree of freedom.

Various training methods exist to modify the weights in order to make the mapping representative of the process data. The gradient descent method is a common one, however quicker methods such as conjugate gradients are available.

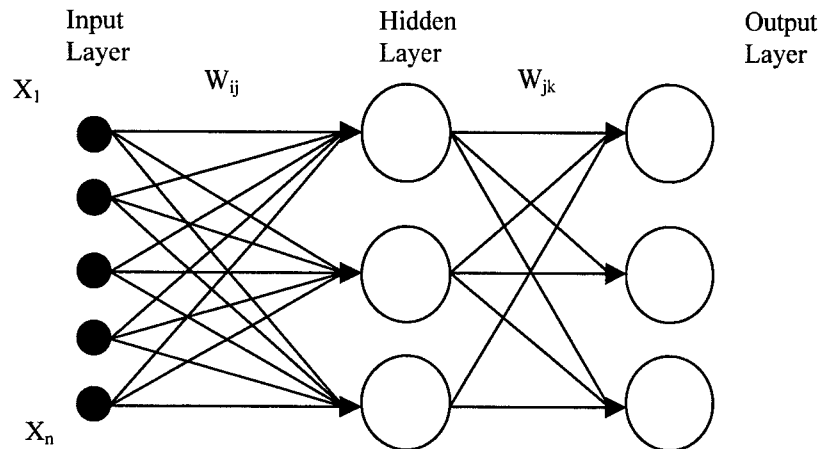


Fig. 1. Multi-layer perceptron structure. N.B. Bias weights are not shown

DEVELOPMENT OF A NEURAL MODEL

When developing a model using process data there are several steps that should be undertaken. These stages are commonly; problem familiarisation, data collection, data preprocessing, training the model, testing the model and commissioning the model. This paper is primarily concerned with 'data cleaning', which mainly occurs at the data preprocessing and training stages, however first we will consider what role the problem familiarisation stage can play in this analysis.

Problem Familiarisation and Data Collection

When embarking on a modelling problem, it is advisable to become familiar with the process. Familiarisation of the real problem can help the modeler to make intuitive decisions throughout the process, but more importantly, it can bring about doubts when a procedure which will not benefit the model is being attempted.

Basic topics such as data availability, format, distribution, variable designation (input or output), variable importance, trends, ranges and units compatibility must be broached. This helps to obtain the most use from the available data, and it helps to judge how many useful data points are available for model construction.

Codes which relate to the production of the product are also important, for example when considering a batch process, each data point will often carry a code relating to its manufacture (a batch number) which may also contain a code relating to the location of its production. Whilst appearing abstract at first, this information may prove useful during outlier detection and correction.

Data Cleaning

Several methods have been used to locate outlying data points and correct their values when modelling data. We will first discuss methods of detection, then correction methods along with missing data treatment and later, we will demonstrate their effectiveness.

Basic Outlier Detection

Once the operating levels of the inputs and outputs have been defined, an obvious check is whether the data set contains any points that violate these limits. Qin et al. [5] used a similar method. It is possible to use ideas relating to the simple physics of a problem to check the validity of data points. Through the familiarisation stage it may be possible to find ways to infer new variables from the data. For example when dealing with steel data, the ratio of two variables is sometimes used as a guide to microstructure. If this ratio exceeds 1, then via an inferred value an outlying data point has been found. However, by looking at the two variables separately it might not have been possible to automatically infer an outlying data point.

A graphical inspection method can also be effective in detecting outliers -- visual correlation checking. For example, if one variable is expected to relate to another in a manner roughly known, outlying data points can be found by identifying those points which do not conform to the general trend of the rest of the group.

Structured Outlier Detection

It is often found in a steel process that certain types of steel are made quite frequently. The result is that there are several examples of a given set of input variables, which should relate to similar output values. Moreover, in some cases there may be data points that come from the same parent cast, undergoing their treatment in either similar or different batches depending upon the quantity produced. These features can be advantageous and a check can be made through the data set to find similar input patterns together with their respective output patterns. By choosing various criteria for similarity, and through knowledge of the system, it is possible to infer an approximate variance in the results for a given set of similar input patterns. Thus, groups that contain points that are outside this range can then be examined further.

If there are many similar examples for a given input pattern, we can classify an outlying data point as the one furthest from the median value. Other information within the data set may also be useful such as information taken at the time of testing. Often in processes such as heat treatment a variable will be assigned to define if a batch passed or failed a specification test. If an outlying point fails to meet a specification due to a process error then it must be excluded, however it is also possible to include data points when the treatment prescription was at fault as these represent valid data and can expand the data set.

Learned Outlier Detection

When a neural model has been constructed, the residuals (the predicted-known values) are often examined. Outlying data points have a feature in which they tend to have large individual training and test errors [5]. However this can be advantageous when detecting outliers, as there are two reasons for high residuals. The first is a data point which the model cannot cope with but one which is correct, and the second is a point which is incorrect, and therefore does not fit in with the model.

One should be aware that when a training data point is incorrect and lies in an area of low data density, with little else to weigh against it, the model will tend to fit toward the incorrect data point; a method to avoid this will be discussed later, in which, the model can detect outlying data points directly.

One method used in this project was to find residuals for both model and test sets that exceed two standard deviations of the mean of the modulus of residuals. With large data sets, this can be limited to the top x number of points. These points are checked for validity by an expert, namely a metallurgist. As shown in Figure 2, the approach can be used in an iterative manner. Points found not to be faulty, but which have high residuals indicate poor areas of the model where additional data may be actively introduced.

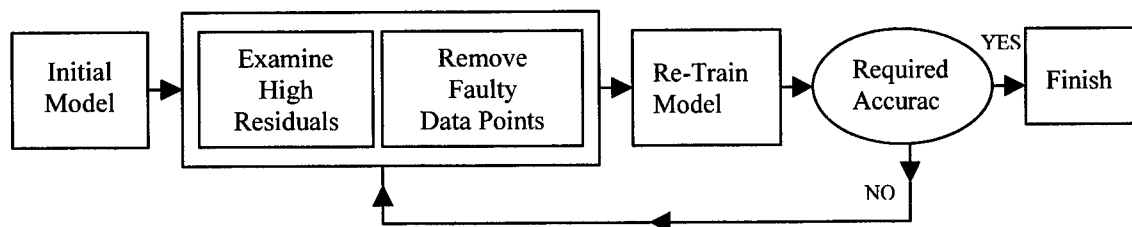


Fig2. The iterative approach of learned outlier detection.

Multivariable Detection

Both Principal Component Analysis (PCA) and Partial Least Squares (PLS) methods have been used for outlier detection [6, 7]. PCA is a method where data with a large number of correlated variables is reduced to a smaller number of principal components, allowing the underlying process to be described. Outlier location is performed on a basis similar to that used for the more-common fault detection [8,9]. Data in certain areas of the input space is clustered when represented in 2 or 3 dimensions, for the first two or three principal components. By investigating points at the edges of these clusters, outliers can be found.

Data Point Correction

When a data point is found to be incorrect, there are several options. First, the data point can be deleted from the set. This is acceptable if other similar examples are available and only a minor effect in the data numbers results. If data are in short supply for that area of the input space, a replacement may be necessary, treating the corrupted value as a missing value. The ideal replacement is that provided by an expert. This can sometimes be achieved if a "partial" expert has access to data from the original time of manufacture.

Rule-based correction can be possible when there are large numbers of outliers to be corrected. This is usually based upon the method of detection, for example when a 'sames' check has been conducted, a specific set of procedures can be employed based on the expert's knowledge for correcting the data points.

If a faulty data point is treated as a missing value and a replacement is required without expert knowledge, there are a number of options. First, the points must be missing at random, i.e., the outliers must occur at random. Missing values can then be found by interpolation. This assumes that the data is uniformly distributed. Linear or more advanced regression can be used, however this is only done with less than 3 consecutive missing points. In the situation in this work, the missing values may be in the input or output space, depending on where the faulty point is, and these techniques can be applied to either case.

PCA and PLS techniques can also be used to treat missing values as stated in Qin et al., 1993 [5], since these techniques basically allow blanks in the data. Therefore, PCA can be applied to the training data with some missing values, then a neural model based on principle components can be constructed. A further approach is to estimate the missing value for an n-dimensional input vector, from knowledge of the other n-1 input variables. The nearest neighbour within the training data set can be identified and the value for the nearest neighbour can be substituted as the missing value. Finally, auto-associative neural networks [10] can be used to fill in the missing value, although this will require a high level of redundancy within the neural network. Sharpe and Solly [11] provide additional information to predict missing data with neural networks.

Outlier Resistant Methods of Training

A further option when dealing with data containing faulty points is to use a network that is resistant to outliers. The sum of squares error function receives the largest contributions from points that have the largest errors. Commonly, outliers have large errors since they are at the extremes of a distribution. Qin [5] mentions the use of a training error to treat large and small errors linearly.

When deciding upon the type of network to use for a model, it is important to consider 'same' outliers. When training a network with ambiguous output values for a fixed input vector, an RBF will not converge if the minimum error function decided upon is greater than the variation of the ambiguity. A MLP however will effectively tend towards the average value available.

CASE STUDY

The process for the model is constructed is described in the background section. The database used to construct the model was collected from five heat treatment sites and a lab testing facility.

Data Collected

The data collected contains information not only about steel composition, size and heat treatment regime, but also data about how it was tested, which batch the steels were treated in, which cast the particular batch stemmed from and whether the required specification set was met. The cleaned data consists of 5711 data points. 194 faulty data points had been identified in the original data.

Table 1. Statistics of cleaned data used in the final model

Variable Name	Type	Min.	Max.	Mean	Std. Dev.
Test Depth	Input	4	140	16.08	9.35
Bar Size	Input	8	381	156.4	83.95
Treatment Site	Input	Binary codes represent 6 locations			
C	Input	0.12	0.63	0.39	0.06
Si	Input	0.11	1.87	0.26	0.04
Mn	Input	0.35	1.75	0.76	0.22
Cr	Input	0.05	3.46	1.04	0.45
Mo	Input	0.01	1.0	0.26	0.14
Ni	Input	0.02	4.21	0.79	0.86
Al	Input	0.005	1.08	0.04	0.09
V	Input	0.001	0.27	0.008	0.023
Temperature at Hardening Stage	Input	820	980	856.9	16.9
Type of Quench at Hardening Stage	Input	Binary codes represent 3 quenches; oil water or air.			
Temperature at Tempering Stage	Input	20	730	604.9	70.7
Ultimate Tensile Strength	Output	516.2	1841	929.1	156.1

Site code and type of quench at the hardening stage are dummy variables which relate to the six locations or the three quench types respectively in gray code binary format. A test reference code, heat treatment batch number, pass/fail statistic, composition code and cast number are used in the outlier detection process.

Data Cleaning

The data cleaning process was followed as described previously, with the following features specific to the process. When structured checking was performed, similar batches were separated into groups of zero differences (duplicate entries), small differences (below a specified process variation of 40 N/mm²), and large differences (in excess of 40 N/mm²). For the small difference and zero difference groups, a median value for all batches within the same analysis with the same heat treatment batch number was used. This prevents the prior probabilities of the data being biased towards the frequency of measurements made on the same batch of steel. Batches with different heat treatment numbers were retained. When dealing with the large difference group, each group of similar batches was investigated separately, as it was expected that either a process or transcription error may have occurred. An automatic rule-based method was devised using expert knowledge to deal with a range of situations where outliers could occur; any remaining points were referred to the expert for correction or deletion.

Experiment

In order to show the effectiveness of the data cleaning process, the following experiment was devised, thus demonstrating a predictive model's generalisation on an unseen test set both before and after cleaning.

When training and evaluating neural network performance, data are partitioned into equal training, validation and test sets [12]. The validation set is used to prevent over-fitting of the training data, by stopping further training when the validation set error starts to rise. Having performed the data cleaning, a test set was randomly selected from the data that had not been deleted by the cleaning process. These points were removed from the cleaned and uncleaned data sets, the cleaned test set was reserved for testing. The remaining data in the cleaned and uncleaned sets formed the training and validation sets. Due to deletions in the cleaned data, points were randomly deleted from the uncleaned set to make the number of data points equal so as not to bias the standard deviation calculation used in the results. The training and validation sets were constructed five times to show the effect of certain data points falling either in the training or validation sets. Training was performed with a MLP containing 6 hidden layer neurons, using back error propagation with gradient descent and momentum. The weights were initialised randomly each time the network was trained.

The five models resulting from the cleaned and uncleaned data were then used to predict ultimate tensile strength values on the unseen test set. With a relatively low proportion of faulty points in the original data, it can be seen from the results in Fig. 3 and Fig. 4, that the cleaned data have a lower standard deviation of residuals for both the training and test sets. Additionally, it can be seen that the predictive accuracy of the cleaned data is more stable when the random ordering of the training and validation sets is changed.

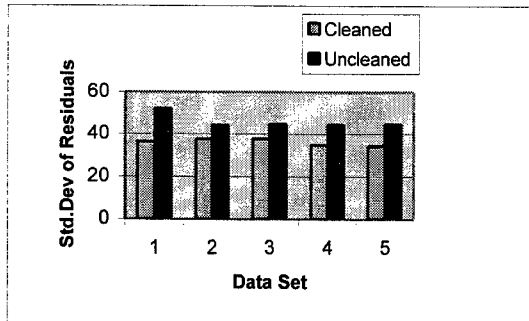


Fig. 3. Training data performance

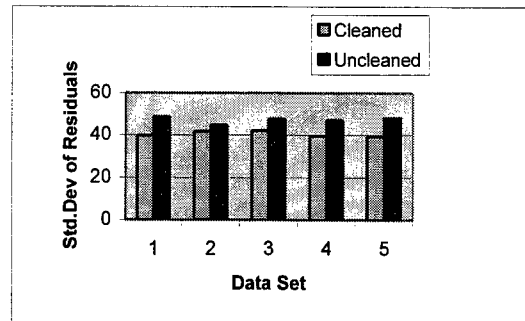


Fig. 4. Testing data performance

CONCLUSION

A variety of data cleaning methods have been presented, which can be applied at the pre-processing and training stages of an empirical model development. The importance of process familiarisation together with expert knowledge has been demonstrated when applying these methodologies to a predictive model for a heat treatment process. Further results will be given at the conference relating to prediction accuracy and generalisation for a range of materials properties. The authors acknowledge funding support from the Materials Forum of companies based within the Sheffield area of the U.K.

REFERENCES

1. T. Cool, H.K.D.H. Bhadeshia, D.J.C. MacKay, 1997. The Yield and Ultimate Tensile Strength of Steel Welds. *Materials Science and Engineering*, 223, 186-200.
2. J. Jones, D.J.C. MacKay, 1996. Neural Network Modelling of the Mechanical Properties of Nickel Base Superalloys. *Proceedings of the Eighth International Symposium*, 417-24.
3. A. Bulsari, E. Hocksell, 1996. Neural Network Systems for Hardened Components. *Steel Technology International*, 133-138.
4. W.S. McCulloch, Pitts, 1943. A Logical Calculus of the Ideas Immanent in Nervous Activity. *Bulletin of Mathematical Biophysics*, 5, 115-133.
5. S. Joe Qin, B. Rajagopal, 1993. Combining Statistics and Expert Systems with Neural Networks for Empirical Process Modelling. *Advances in Instrumentation and Control*, 48,3, 1711-1720.
6. C. Dumortier, P. Leher, P. Krupa, A. Charlier, 1998. Statistical Modelling of Mechanical Properties of Micro-alloyed Steels by Application of Artificial Neural Networks. *Mat.Sci. Forum*, 284-286, 393-400.
7. P. Nomikos, J.F. MacGregor, 1994. Monitoring Batch Processes Using Multiway Principal Component Analysis. *AIChE Journal*, 40,8 1361-1375.
8. C. Thomas, T. Wada, D.E. Seborg, 1996. Principal Component Analysis Applied to Process Monitoring of an Industrial Distillation Column. 1996 IFAC 13th Triennial World Congress, 61-66.
9. N. Jaleel, M. Fiocco, J.R. Leigh, 1994. Monitoring Estimation and Predictive Control Based on Statistical Techniques. *Proceedings of the American Control Conference*, 328-329.
10. M.A. Kramer, 1992. Autoassociative Neural Networks. *Computers Chem. Engng*, 16,4, 313-328.
11. P.K. Sharpe, R.J. Solly, 1995. Dealing with missing values in neural network-based diagnostic systems. *Neural Computing and Applications*, 3,73-77.
12. L. Tarassenko, 1998. A Guide to Neural Computing Applications. *Neural Comput. Appl. Forum*, 74-75.

Intelligent Measurement System Confirmation

P. H. Osanna, M.N. Durakbasa

Department for Interchangeable Manufacturing and Industrial Metrology
Vienna University of Technology, A-1040 Wien, Austria

ABSTRACT

Through the dynamification of testing intervals the flexibility level of confirmation systems for measuring equipment can be increased considerably. Since confirmation of inspection, measuring and testing equipment is a significant part of quality management and an essential requirement for the entire production process --especially with the increasing demands of micro- and nanotechnology-- it is absolutely necessary to increase the efficiency of confirmation systems. Through a special method developed for this purpose, the flexibility level and efficiency of confirmation systems can be achieved and the expenses can be substantially reduced by the use of fuzzy logic for dynamification of testing intervals.

INTRODUCTION

The permanent increasing of quality standards, world wide competition, as well as the legislation of regulation of the product responsibility, require not only a proper documentation of the measurement data of the production, but also the continuous supervision of measuring and test equipment. Especially in modern computer-integrated production, testing devices are often connected directly with the manufacturing process. This causes direct or indirect influences on the quality level, therefore the confirmation and management of measuring and test equipment is becoming a significant part of the quality management for the entire production [1, 2].

The confirmation of measuring equipment is an essential quality requirement for modern production especially at the higher demands of micro and nanotechnology. The efficiency of the confirmation can be increased and expenses can be reduced substantially through computer assistance with flexible testing intervals. For this purpose a special method has been developed at the Department for Interchangeable Manufacturing and Industrial Metrology, with which an increase of the flexibility level and efficiency of a system for the intelligent management and confirmation of inspection, measuring and test equipment can be achieved.

CONFIRMATION AND MANAGEMENT OF MEASURING INSTRUMENTS

Measuring and test equipment connected with the production flow, are subject to a relatively high wear through constant use. Therefore a universal function capability for the entire production flow is absolutely necessary. This confirms the necessity of regular examination and documentation of these measuring devices. But it is possible that unused measuring and test equipment lose their fitness for use through physical or chemical influences. Also new measuring and test equipment may have errors and they are subject to wear at later stages of use. The management of measuring and test equipment should ensure that the used measuring equipment function at all times. Furthermore the operational precision of the used measuring and test equipment or their maintenance and repair should be ensured by management of measuring equipment [3]. Certain systems of confirmation of measuring equipment are already employed as criterion for order replacement. Many enterprises, especially automobile production, electronic and steel industries, demand confirmation and management system of measuring equipment from their suppliers [4].

The system of management of measuring equipment requires the labeling of the measuring equipment, so that a clear identification can be guaranteed. In computer aided control and checking of measuring equipment basic data are provided to each measuring device, composed essentially of:

- identification data (group number, designation, producer, location of usage),
- descriptive data (unit of measurement, range of measurement, resolution),

- location of use, location of deposit, date of the lending / return, user,
- status data (status of testing equipment),
- test monitoring (test standard, last / neighbour test, testing interval, testing factor) and
- informative data (temperature range, accessory, maintenance costs, purchase date and price).

Measuring equipment management systems also can assist in the lending of the measuring equipment from the store. Furthermore each measuring equipment will be marked according to its whereabouts such as:

- blocked,
- in operation,
- in the store,
- under repair,
- scrap,
- not traceable.

As a first step, checking plans for groups of measuring instruments shall be prepared. These can be associated with the basic data of the individual measuring equipment. Checking plans for certain groups of measuring equipment should be stored under the respective group number, so that they can apply to all measuring equipment of this group.

The intelligent confirmation and management of measuring equipment provide a specific record of checking results to each measuring device. For each measuring equipment recorded in the data base control charts displaying the results and the judgement of the finally accomplished test shall be laid out.. At variable features values appear in their nominal value with relevant specification of tolerance. The actual deviations are to be entered with a corresponding sign. A more automatic variance comparison produces a finding in the form "ok" or "n.ok". The automatically prepared finding should be changeable by the user. At attributive features only an entry of the form "ok" or "n.ok" appears.

The evaluation function of the system should give an overview of the individual history of the measuring equipment. This allows the observation and analysis of the development of a measuring device according to its characteristics. The results of the last test as well as the entered deviations should be indicated. At variable features the measuring behavior can be presented additionally in the form of a graph.

DETERMINING OF FLEXIBLE TESTING INTERVALS WITH COMPUTER-ASSISTANCE

A very sensitive criterion for the continuous monitoring of the measuring instruments is the checking interval. Checking interval is the distance between two tests following each other. It is also important to determine the confirmation intervals of measuring equipment, because all relevant influence quantities, e.g. intended use or claim, should be taken into account. In general the checking intervals can be determined as time interval, as limit of the number of the uses or as combination of both. All measuring and test equipment according to the scope are to be examined in determined distances. The examinations are to be determined in virtue of type, stability, intended purpose and utilization frequency of the measuring and test equipment.

On the basis of results of preceding calibrations the intervals are to be shortened or lengthened to secure the continuous precision. To get an adequately small uncertainty of measurement, the confirmation intervals should be chosen as short as possible but, on the other hand, this will lend itself to a high rate of equipment utilization. The system must ensure, that the measuring and test equipment will be calibrated according to the determined timetable. If the checking interval is exceeded, the measuring equipment has to be marked and blocked for further application.

If instruments are used without any practical experience the interval can be determined only approximately. In that case the experiences of the similar instruments as well as the different factors appearing at the measurement elements should be considered. In case of doubt, the interval is to be set shorter than anticipated and to be corrected in the next examination if necessary [5].

"Optimal Interval" is one where the total costs are a minimum. If the interval is chosen too small, the checking costs go up, because there is more checking in the equal period than necessary. If the interval is chosen too large, there rises the probability to find the measuring and test equipment as "inadmissible" at the next checking. This means, that longer time would be used with an inadmissible measuring equipment and thus the quality of the delivered products would be low. Poor quality always implies an increase of costs.

The size of the optimal interval depends on a number of different elements: frequency of utilization, mode of using, behavior of abrasion, consequences at lapses, permissible tolerance range, number of the users, status in the calibration chain, etc. Since these elements are not temporal constant, the optimal interval can not be a constant size.

DYNAMIFICATION OF THE INTERVAL BY USING FUZZY LOGIC

Through the dynamification of the checking interval it is possible to increase the level of flexibility of the measuring equipment management system. Refraining from fixed checking intervals naturally means an increased expense for fixing the intervals. However the dynamification allows to optimize the measuring expenditure and to increase the reliability of the measuring and test equipment. The construction of the fuzzy system for the evaluation of a measuring instrument is shown in Figure 1.

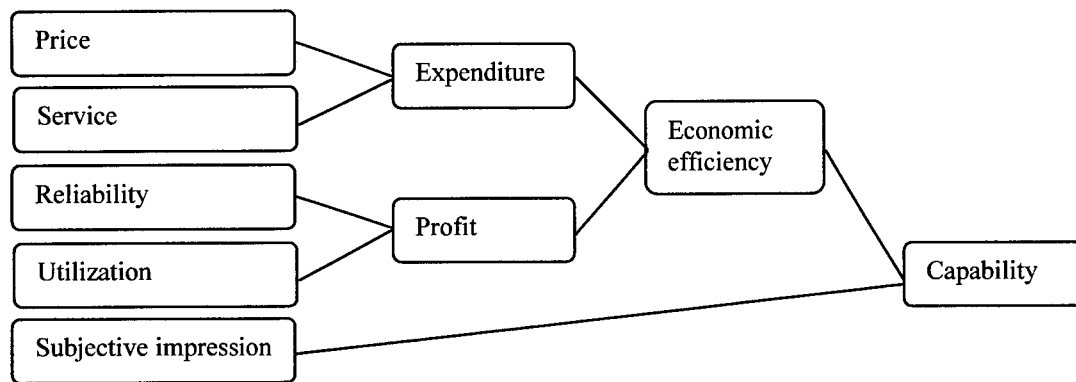


Fig. 1. Fuzzy system construction for measuring instrument evaluation.

For this purpose there is a method for the determination of the testing interval based on expert knowledge, which has been developed at the Department for Interchangeable Manufacturing and Industrial Metrology. This method makes the estimate and the insertion of the experiences unnecessary, since all "estimate"-processes are accomplished using computer aid and practical knowledge (experience from earlier tests, experiences from the testing, etc.) is implemented in mathematical algorithms. This method is based on demand of the environmental conditions of the past as well as of the expected future. This happens exclusively through the application of fuzzy logic, a polyvalent logic, which operates mainly with linguistic variables [6].

Central characteristic within the computer aided interval dynamification is the so-called test-characteristic. This characteristic describes the condition of the testing equipment at time of testing. The test-characteristic is a result of the considered testing. The range of results of the test-characteristic is the interval in number between 0 (inadmissible condition of measuring device) and 1 (optimum condition).

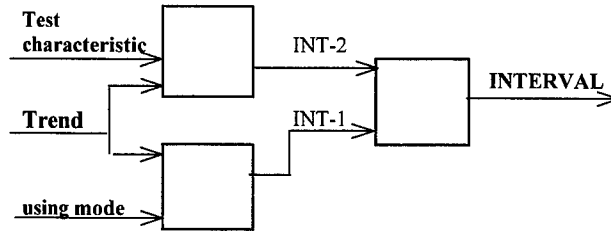


Fig. 2. Substitute model.

The sizes necessary for the dynamification of the interval are shown in Figure 2. Therein one recognizes, that beside the test-characteristic two further input sizes: the trend which is determined again by the test-characteristic, and the mode of use. This trend indicates the temporal course of the test-characteristic-curve. The trend will be at constant environmental conditions a linear slope. The mode of use implies the behavior of abrasion, frequency of utilization, pollution degree, operation temperature, the storage at not on the operation etc. This becomes thereby exclusively the expected mode of using until the next checking.

Since there are more than two input sizes, reference shall be made to a substitute model (Figure 2), which consists of two subsystems with only two input sizes. The testing interval is given as an output size in the range of 3 and 24, which indicates the number of months. For establishing the fuzzy sets INT-1 and INT-2 the five linguistic terms are divided up according to the standard form. In the fuzzy set of testing interval, the INT-1 and INT-2 are other combined with each other to produce the optimal interval (Figure 3).

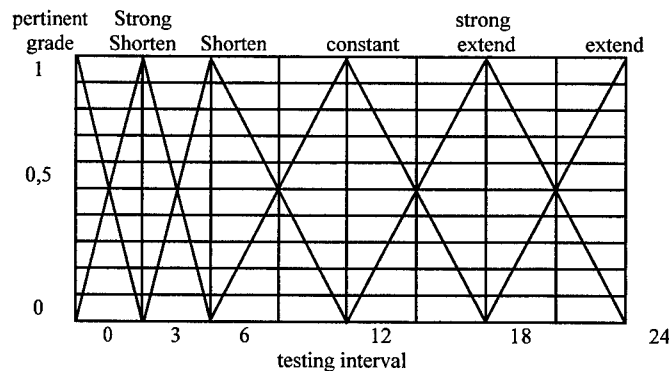


Fig. 3. Fuzzy sets of the testing intervals.

The control basis is created by an iterative method. Figure 4 shows the control basis. All sizes are in complex interaction with each other. The method of fuzzy logic provides for a reproducible system to find the optimal interval. Requirement is that the internal fuzzy algorithms are formed in virtue of reliable expert knowledge.

DISCUSSION AND CONCLUDING REMARKS

The confirmation and management system of measuring equipment should ensure that the measuring instruments function at all times. Through computer aided systems the expense of documentation because of the multitude of use of measuring equipment, as well as the expenses concerning the establishment of the checking plans are lowered. These systems should provide for complete and orderly checking and for the traceability of measuring and calibration results.

By means of fuzzy logic it is possible to create a system which allows the quantitative evaluation of the quality of measuring equipment. Because of many relevant influence quantities, this system can be used for optimizing the confirmation intervals. Further, through the use of this systems within the management of manufacturing equipment the checking expense climbs only conditionally, while the profit gains increase in

greater proportions and the estimate of the testing intervals based on subjective impressions are objectified and made reproducible. The economic benefits and the quality of the confirmation of measuring and test equipment will increase substantially.

		INT-1				
		TREND DEVIATION				
		SS	S	L	R	SR
application mode	SG	E	E	SE	SE	SE
	PM	SH	C	E	E	SE
	M	SH	SH	C	E	E
	PR	SSH	SH	SH	C	E
	PV	SSH	SSH	SH	SH	C
		→ INT-1				

SG	storage
PM	Precision measur.room
M	measuring room
PR	Production, rough
PV	Production, very rough

SS	steep slop
S	Slop
L	Linear
R	Rising
SR	Steep rising

		INT-2				
		TREND DEVIATION				
		SS	S	L	R	SR
test character	VP	SSH	SSH	SSH	SSH	SSH
	P	SSH	SSH	SH	SH	E
	M	SSH	S	C	E	SE
	G	SH	C	E	SE	SE
	VG	C	E	SE	SE	SE
		→ INT-2				

VP	Very poor
P	Poor
M	Middle
G	Gut
VG	Very gut

SSH	Strong shorten
SH	Shorten
C	Constant
E	Extend
SE	Strong extend

		INT-2				
		SSH	SH	C	E	SE
INT-1	SS	SSH	SSH	SSH	SSH	SH
	S	SSH	SSH	SH	SH	C
	C	SH	SH	C	C	E
	E	SH	C	C	E	SE
	SE	C	E	E	SE	SE
		→ TESTING -INTERVAL				

SSH	Strong shorten
SH	Shorten
C	Constant
E	Extend
SE	Strong extend

Fig. 4. Control basis.

REFERENCES

1. EN/ISO 9000-1, 1994. Quality Management and Quality Assurance Standards - Part 1: Guidelines for Selection and Use.
2. EN/ISO 9004-1, 1994. Quality Management and Quality Systems Elements; Guidelines.
3. Durakbasa, M.N., 1996. Ueberwachung, Verwaltung und Faehigkeitsuntersuchungen von Pruefmitteln, 113(4), 286-287.
4. ON M 1340, 1991. Pruefmittelueberwachung. Austrian Standard.
5. Durakbasa, M.N., Pfeiffermann, G.G., 1997. Computer Aided Confirmation and Management of Inspection, Measuring and Test Equipment in the Quality Management Systems, Proceedings "Measurement 97", Inter. Conf. on Measurement, Smolenice, Slovak Republic, May, 63-66.
6. Mamdani, E., Sedrak, A., 1995. An Experiment in Linguistic Synthesis with a Fuzzy Logic Controller; International Journal of Man-Machine-Studies, 7, 1-13.

Simulation of Dynamic Properties of Nuclear Meters Applied in Coal Preparation Control Systems

S. Cierpisz

Silesian Technical University
Department of EE and Process Control in Mining
40-100 Gliwice ul. Akademicka 2, Poland
Fax: (+48)32 2371537 Email: cierpisz@zeus.polsl.gliwice.pl.

ABSTRACT

On-line nuclear meters have been in use in the coal industry for many years. They have been utilized for coal quality monitoring, in control systems for coal blending or for treating coal in the heavy media separation process. Their operation is based on the scattering or absorption of incident gamma radiation, and the derived density or ash value is the result of a time-averaged measurement. In this paper, dynamic models of ash monitors have been presented and discussed. The analysis of monitors with constant time of measurement shows that it is possible to determine optimal time for which the dynamic error is the smallest. The analysis also shows that the best results are given by the monitors in which the time of measurement is variable and adapts to changes of the input signal. An example of a coal blending control system which stabilizes coal quality of the blend has been analyzed. The simulation of the system operation has been performed with the use of the Matlab (Simulink) program package.

INTRODUCTION

On-line nuclear meters such as heavy media and coal slurry densitometers or ash content in coals monitors have been in use in the coal industry for many years. They have been utilized for coal quality monitoring, in the control systems for coal blending, or for separating coals in the heavy media separation process. Their operation is based on the scattering or the absorption of incident gamma radiation, and the derived density or ash value is the result of a time-averaged measurement over a period of tens of seconds or even a few minutes. Such monitors cannot produce an exact dynamic record of the rapid variations in the parameter to be measured.

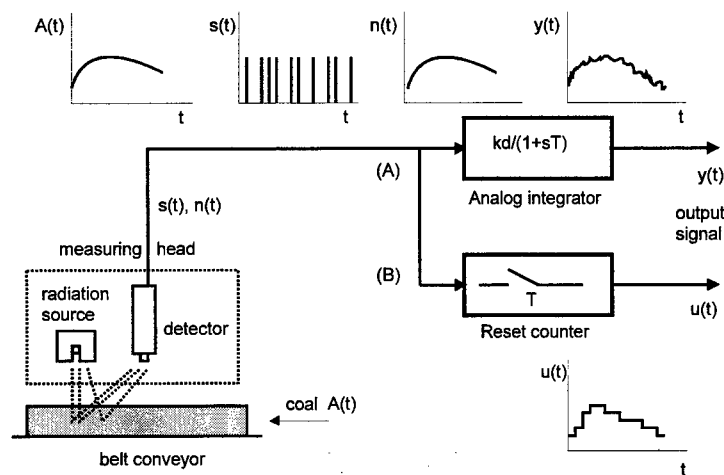


Fig.1. Diagram of on-line ash monitor

A general scheme of an on-line ash monitor is shown in the Figure 1. The electronic circuit with an analog integrator with the time constant T_i for a series of pulses from a detector is shown in the part (A) of the

Figure the part (B). At present in the majority of monitors a digital counter is used. The output signal from the detector (scintillation counter) is always a stochastic signal, regardless of the character of the input signal (e.i. ash content) modulating the intensity of the detected radiation beam. The longer the averaging time the higher the statistical (static) accuracy of the monitor. At the same time, if the input signal varies, the dynamic error of the measurement is higher. This suggests that for a given shape of the input signal and a given structure of the monitor circuit, one can find an optimal averaging time of input pulses which gives the minimum dynamic error according to the accepted criteria. Furthermore, this leads to the application of a circuit with an adapting time constant. If the input signal is, for example, a step function and the ash monitor is to reproduce this change, the time constant should be small at the beginning of the measurement to speed-up the reaction of the meter and then it should become greater to read-out accurately the new value of the ash content.

A concept of an ash monitor with a time constant (or a time of measurement) adapting to variations of an input signal (ash content) has been analyzed. Such a system allows to speed-up the reaction of the instrument to rapid variations of ash content and at the same time to achieve better statistical accuracy for a longer period of time. This is particularly important in closed loop control systems or in splitting of a coal stream to different products.

DYNAMIC PROPERTIES OF ON-LINE ASH MONITOR

A detailed description of the operation of the ash monitor can be found in [3,5,6]. Let us assume that one wishes to measure accurately a step change of ash content $A(t)$ controlled by the monitor with the digital counter of pulses shown in the Figure 2b. That means that the shape of the output signal $u(t)$ should closely resemble the input step function so as to minimize the distance between $A(t)$ and the output signal $u(t)$.

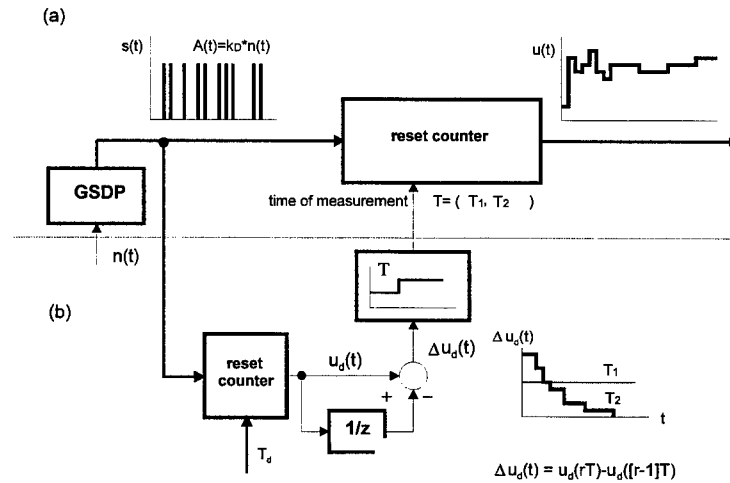


Fig.2. Model of ash monitor with reset counter.

For small changes in ash content $A(t)$, let us assume a linear relation with mean intensity $n(t)$ of pulses $s(t)$:

$$n(t) = n_0 + \Delta n(t) = k_d * (A_0 + \Delta A(t)) = k_d * A(t) \quad 1.$$

The stochastic signal $u(t)$ can be calculated from the equation:

$$u(rT) = \sum_{k=1}^{k=l} N_k \quad 2.$$

where N_k is the number of pulses which appear in the interval of time $t_k = (r-1)T - rT$.

Mean value $M_u(rT)$ and variance $D_u(rT)$ of output signal $u(t)$ from the digital counter in Figure 1b can be calculated [5] from the equations (in time intervals $t_k = (r-1)T - rT$):

$$M_u(rT) = \int_{(r-1)T}^{rT} n(t) dt \quad 3.$$

$$D_u(rT) = \int_{(r-1)T}^{rT} n(t) dt \quad 4.$$

The distance between the input signal $A(t)$ and the output signal $u(t)$ is often defined [5] as the mean-square relation:

$$L = \int_0^{T_0} ([k_d * A(t) - M_u(t)]^2 + D_u(t)) dt \quad 5.$$

Equation 5 can give analytical solution for optimal parameters of the monitor basically for linear systems with constant parameters [5,6]. For nonlinear adaptive circuits the computer simulation analysis can be used for optimization of the monitor operation.

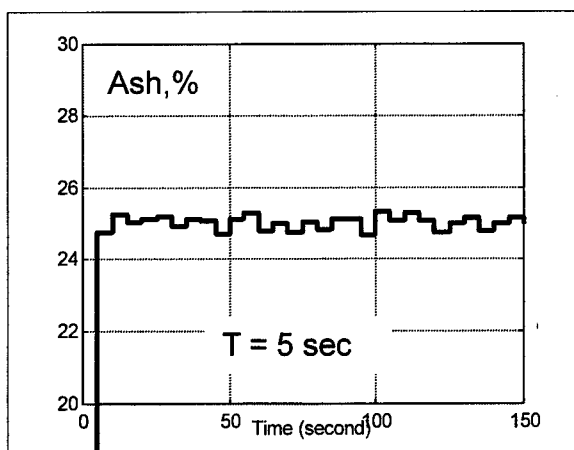


Fig.3. Response of ash monitor to step change of ash content (25%) $T = 5$ s.

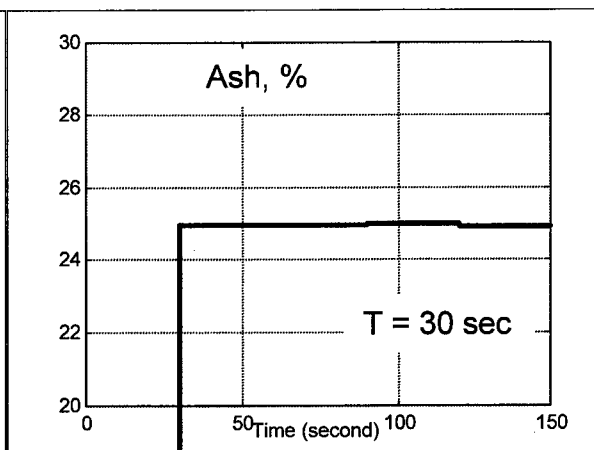


Fig.4. Response of ash monitor to step change of ash content (25%) $T = 30$ s.

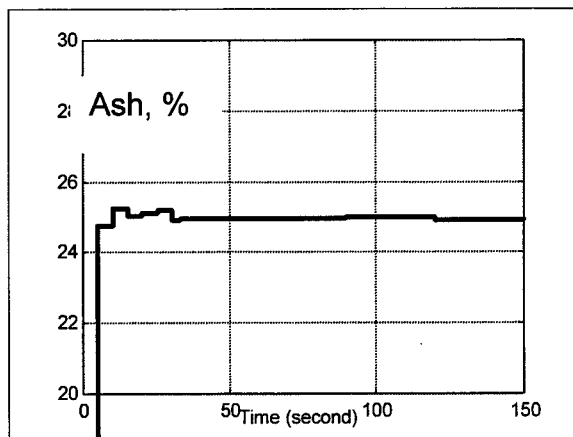


Fig.5. Response of ash monitor to step change of ash content (25%) adaptive $T(5$ and 30 s).

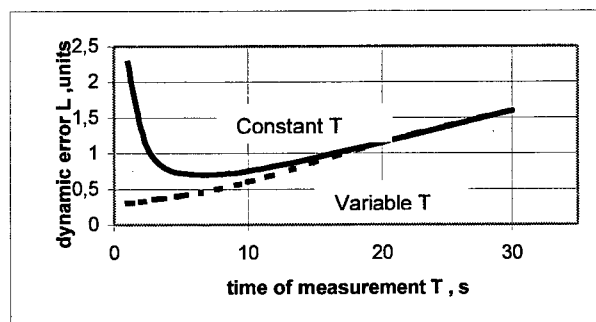


Fig.6. Dynamic error of measurement as a function of the time of measurement T .

The simplified model of an on-line ash monitor with a reset counter is shown in Figure 2 in which the signal $u(t)$ is the number of pulses $s(t)$ which appeared during the intervals $(r-1)T - rT$ and GPDN is a Generator of Poisson Discrete Noise with the controlled mean intensity of pulses $n(t)$. The model of the monitor with reset counter of pulses and a constant time of measurement is shown in Figure 2a. This model can be modified to the circuit with a variable time of measurement adapting to changes of the input signal $n(t)$ shown in the Figure 2b. Changes of ash content with time $u_d(t)$ and their speed (derivative) $\Delta u_d(t)$ are detected in the additional reset counter with a short time of the measurement T_d . The time of the measurement T of the basic counter adjusts to the value of $\Delta u_d(t)$. In the simplest case T can have two values, it can be short (T_1) for high values of $\Delta u_d(t)$ and long (T_2) for low values of $\Delta u_d(t)$ (for instance $T=T_1=5s$, $T=T_2=30s$) as it is shown in the Figure 2. The response of the ash monitor to step changes of ash content in coal has been simulated with the use of MATLAB Simulink software package and has been shown in Figures 3,4,5. The dynamic error of the measurement or the distance between input and output signals has been calculated in this case from the equation:

$$L_u = \int_0^{T_0} [k_d * A(t) - u(rT)]^2 dt \quad 6.$$

The error L_u shown in Figure 6 depends on the time of measurement T ; it is high for both short and long time of measurement and is the smallest for the optimal value of T . The ash monitor with the time of measurement T adapting to changes of ash content gives a better response and a smaller dynamic error L_u in comparison with the conventional type of monitor.

CONTROL SYSTEM

An example of a coal blending control system with an on-line radiometric ash monitor is shown in Figure 7. The blend is produced from a concentrate and raw coal. Proportion of tonnages of these components is set by the controller on the basis of indications of belt weighers and the desired value of ash content in the blend. The proportion of tonnages is continuously corrected by the controller according to the periodical measurements of the ash monitor. A simplified scheme of the control system is shown in Figure 8. The operation of the control system has been simulated with the use of Matlab Simulink software package.

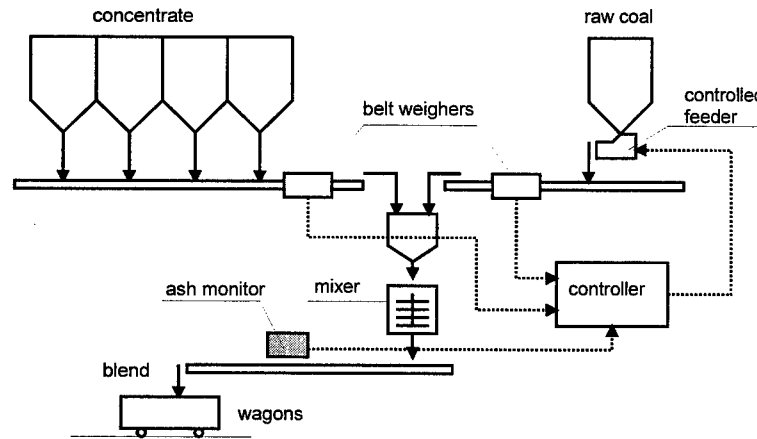


Fig.7. Coal blending control system.

Response of the control system to a step change of the ash content desired value A_{des} is shown in Figures 9 and 10. The change of the ash content in the blend in the Figure 10 corresponds to the ash monitor with constant time of measurement, whereas Figure 9 shows the response of the system with the ash monitor having the variable (adapting) time of measurement ($T=T_1$ or $T=T_2$ for big or small error of the control $e(t)=A_{des}-A(t)$). The application of the adaptive ash monitor gives better results of the control than the system with the constant time of measurement T .

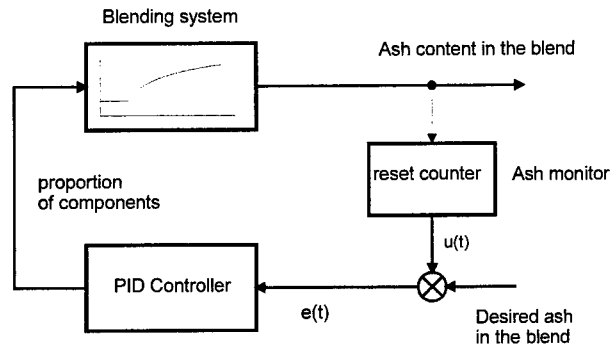


Fig.8. Simplified scheme of the control system.

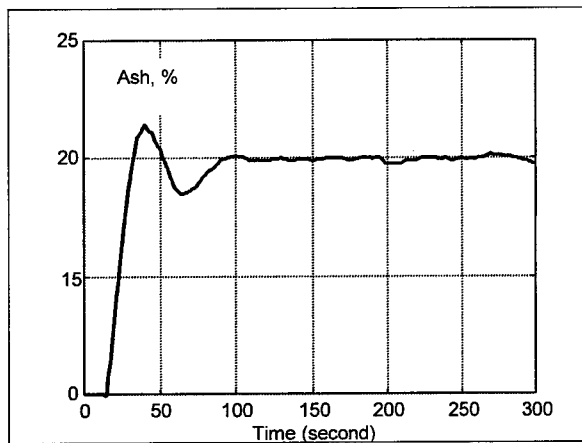


Fig.9. Response of the control system to step change of ash content desired value (adaptive T).

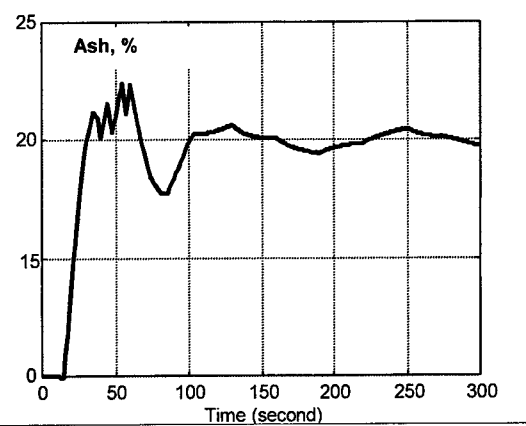


Fig.10. Response of the control system to step change of ash content desired value (constant T).

CONCLUSIONS

One can consider the electronic circuit of an on-line ash monitor to be a generator of Poisson discrete noise with the mean intensity of pulses controlled by the measured signal (ash content). The analysis of the circuit dynamics shows that it is possible to determine an optimal measurement time, which gives minimum dynamic error of the measurement.

Dynamic properties of an on-line ash monitor can be upgraded in circuits in which the time of measurement adapts to changes of the input signal. In the case of an ash monitor used for monitoring, the time of measurement should adapt to the differential input signal and in the case of a closed loop system, the time of measurement should generally adapt to the magnitude of the error of the control.

REFERENCES

1. B.D. Sowerby, J.S. Watt, 1990. Development of nuclear techniques for on-line analysis in the coal industry. *Nuclear Instruments and Methods in Physics Research*; A299(1-3);642-647.
2. B.D. Sowerby, 1991. Nuclear techniques in the coal industry. in: *Inter. Symp. IAEA, Vienna*, 3-31.
3. A.T. Kirchner, 1991. On-line analysis of coal. *IEACR/40*.
4. S. Cierpisz, T. Inoue, 1994. *Process Control and Monitoring in Coal Preparation*. 12th ICPC, Krakow.
5. S. Cierpisz, T. Sikora, 1994. Coal Quality Monitoring and Control in Poland. *J. Coal Qual.*, 13(1), 30-33.
6. S. Cierpisz, T. Sikora, 1994. Dynamic Properties of on-line Ash Monitors. *J. Coal Qual.*, 13(1), 28-30.

ACOUSTIC EMISSIONS MONITORING OF SAG MILL PERFORMANCE

S.J. Spencer, J.J. Campbell, K.R. Weller and Y. Liu

CSIRO Minerals, PO Box 883, Kenmore, QLD 4069, Australia

ABSTRACT

Particle grinding is a key final stage in the comminution process used to liberate minerals from gangue in mineral processing plants. Semi-autogenous grinding (SAG) occurs in tumbling mills which use both large ore particles and steel balls as grinding media. Grinding is the most energy intensive of the mineral processing unit operations and hence the optimisation of the process in terms of throughput rate, mill power consumption, product particle size distribution and mill liner wear is of considerable interest to industry.

Surface vibration (commonly termed acoustic emissions) monitoring is widely used as a non-invasive, low cost means of monitoring normally inaccessible attributes of processes or equipment operation. The current investigation monitors vibrations in an operational SAG mill by an accelerometer attached to the outside of the rotating shell. Attempts are made to interpret the vibrations recorded in terms of the internal state of the SAG mill for a conditional experimental program over a wide range of mill operating conditions.

Results of this study support the view that higher feed rate dynamic steady states correspond to an increased charge mass, with enhanced cushioning of grinding media impacts on the liner due to an increase in the intervening charge volume. An increase in mill rotation speed results in grinding media being lifted higher and more often directly impacting on the shell liner, thus increasing acoustic emissions. Increased pulp density aids in damping collisions that generate acoustic emissions by increasing the resistance to transport of grinding media through the charge. Addition of grinding balls results in more high energy impact events between balls and liner.

The relationships shown in the paper indicate promise for acoustic emissions measurement to be used as part of a system for both process control and condition monitoring of SAG mills.

INTRODUCTION

The final stages of comminution (size reduction) of ore in mineral processing plants for the liberation of valuable minerals from gangue (non-valuable rock) is usually achieved by grinding in tumbling mills. These are rotating drums that contain a charge of loosely packed grinding media (rods, balls or large rocks), occupying less than half the internal volume of the mill, which aid in comminution of particles in a slurry of ore and water. Ore particles are typically reduced in size from about 5-250 mm to 10-300 μm by shatter, cleavage and abrasion events involving cataracting or cascading grinding media and in some instances, contact with the liner and/or charge lifter bars on the inside of the mill shell. Large (up to order ten meters in diameter by five meters in length) SAG (semi-autogenous grinding) mills are favoured in the mineral processing industry for primary grinding of large tonnages (order 100-1000 tonnes per hour feed rate) of many ores [1,2]. Conventional grinding mills typically contain a steel charge of about 35-40 vol.% while SAG mills require a steel charge of only 6-18 vol.%.

Acoustic emissions or surface vibration technologies have been used in various forms to investigate and control the performance of autogenous grinding (AG) mills, ball mills and cyclone classifiers. Control of power draft and mill load in AG and SAG mills has traditionally been by load cells estimating the charge mass. However, acoustic emissions registered by external dual microphone systems have been used to monitor the changing level of impact of rotating charge on an AG mill shell and hence as an input for mill control [3]. The same study found that acoustic emissions are indicators of pulp density and viscosity [3]. Estimation of effective pulp density and viscosity by the magnitude of acoustic emissions has also been

achieved for laboratory batch ball mills [4]. Mill acoustic emissions have also been shown to indicate charge size distribution, ore breakage rates, and ore character in batch ball mills [5,6].

Acoustic emissions monitoring has been used to analyse the performance of a hydrocyclone mineral classification device [7]. Results of this study indicated that acoustic emission frequency domain spectral features in frequencies up to ~50 HZ are sensitive to operating conditions. Relationships were derived between the operating parameters of the cyclone and the spectral and statistical characteristics of the acoustic emissions. Surface vibration monitoring has also been used to study the feed distribution characteristics of parallel Dense Medium (DM) cyclone classification devices in a coal preparation plant [8]. The method is based on the concept that accelerometer sensor monitoring of vibrations at various points on the external surface of a cyclone can yield the frequency and strength of particle impacts on the inside of the hydrocyclone wall. Differences in internal operating characteristics of DM cyclones were detected by surface vibrations as a function of feed conditions and sensor position.

The operation of a SAG mill results in the generation of high frequency surface waves on the outside of the rotating shell due to collision events within the mill. Monitoring of surface vibration waves with an accelerometer attached to the mill shell therefore provides information on events inside the mill, particularly impacts of grinding media on the liner. However, measurements of surface vibrations on the outside of the shell do not simply reflect local impact events on the inside of the liner. All of the components of the mill behave to some extent as elastic media, permitting the propagation of waves generated by collision events within the mill, 'flexing' of the mill shell during rotation and external sources such as the drive motor and girth gear. An accelerometer mounted on the outside of the shell registers normal acceleration due to waves transversely propagating around the shell. These waves are damped in accordance with the properties of the elastic media between the point of wave registration and the origin of the causative event. Hence vibrational events as measured by an accelerometer can be expected to be due to causative events over a limited range of locations within the mill and associated assembly. However, for the preliminary analysis reported in this paper it is assumed that the vibrations are locally generated by collision events inside the mill, adjacent to the accelerometer.

A non-intrusive means of quantifying the spatial position and intensity of the various types of grinding behaviour in SAG mills would be very useful for process monitoring and control. SAG mill operators are particularly keen to utilise a technique that provides a reliable measure of mill load as part of mill control strategy. Monitoring of the frequency and energy of grinding media direct impacts on the shell above the charge region would also be very useful for control of SAG mill liner wear.

APPARATUS

The surface vibration monitoring apparatus to be mounted on the mill shell consists of an accelerometer connected to a charge amplifier. The output from the charge amplifier is connected to a microphone belt-pack transmitter powered by a rechargeable battery connected to two solar panels mounted on opposite sides of the mill. Transmitted data is received using a microphone wireless receiver with two modified extended antennae. Receiver output is connected to a terminal block which is linked to a laptop computer by a fast data acquisition card.

A magnetic proximity pad was mounted on the mill at the 3 o'clock position looking from the discharge end. The detector/switch was mounted off the mill and connected to the terminal block mentioned above. The switching signal from the proximity detector is used as a trigger for logging of the accelerometer signal.

The data acquisition software (written in LabView) can be triggered manually or digitally. Triggering occurs when the magnetic pad and detector/switch comes in close proximity as the mill rotates. Data acquisition then begins with the data being read into a rolling buffer at a rate adjustable up to 100k samples/s. Unprocessed raw data was saved in binary format for further data processing analysis.

EXPERIMENTAL DESIGN

A series of surface vibrations monitoring runs based on conditional experimental design were conducted on the SAG mill at the Red Dome gold mine in Australia. A total of 23 test runs were conducted to investigate how surface vibration features changed with one manipulated operating variable at a time at dynamic steady state (steady power draw) conditions. The manipulated operating variables were tonnage feed rate, mill speed, mill discharge density and ball addition. In addition to acquiring surface vibration information, both control system data and physical plant measurements were taken at each set of conditions to confirm test run validity. Data referred to in this paper were acquired with the acquisition speed set at 5×10^4 scans/s for a duration of 1×10^6 scans. The tonnage feed rate was 170-200 tph (mill speed - 11.8 rpm, no ball addition, pulp density - 72% solids w/w). The mill rotation speed range was 12.3-13.8 rpm (tonnage - 210 tph, no ball addition, pulp density - 72% solids w/w). A number of tests were conducted at intermediate conditions for the tonnage and speed. Experiments were also performed with grinding ball addition and pulp density.

SIGNAL PROCESSING TECHNIQUES

The goal of the data analysis techniques is to derive quantitative measures and qualitative visualisations based on the response of the accelerometer to shell vibrations that can be correlated with SAG mill operating conditions. Vibration measurements may then be used for process condition monitoring and as an input to unit control. The measures may also be useful for inference of the rate of liner wear as a function of operating conditions in with the SAG mill.

Surface vibrational waves as registered by an accelerometer are characterised by a wide variety of measures. The first step in data processing prior to deriving any of these measures is to truncate the data to an integral number of mill rotation periods. This is done in order to ensure that there is no bias in the data due to the sensor detecting changes in mill conditions as a function of rotational position of the outer shell.

The concept of a shell surface vibration *event* is important in the data processing. Such an event is defined as a positive deviation from nil accelerometer response. The amplitude is taken as the peak accelerometer response associated with a positive acceleration. This is in accordance with a propagating surface wave inducing a positive acceleration in an accelerometer corresponding to a normal stress outwards from the shell. It is hypothesised that collision events within the mill, particularly grinding media/liner events, will induce a strain that will propagate as a wave to the outside of the shell and be initially sensed as a positive acceleration. A wave train due to a collisional event should be composed of an initial relatively large, positive acceleration followed by negative and positive oscillations of rapidly decreasing amplitude. It is expected that the accelerometer will detect only the first few oscillations of any wave train associated with a particular collisional event. Negative accelerations are interpreted as part of a wave train belonging to a previous positive acceleration and are hence discarded in terms of registering distinct events. Subsequent positive oscillations in a wave train are expected to be highly damped due to the low elasticity and high damping properties of the liner and the outer shell. Hence it is reasonable that each sequence of positive acceleration defines a vibrational event caused by a particular media/media or media/liner collision within the mill.

The surface vibration measures related to the sampled signal that are used to characterise surface vibrational waves are as follows. Mean, standard deviation, Power Spectral Density (PSD) plots, amplitude histograms and the ratio of numbers of large to small amplitude samples. Power spectral density derived from the amplitude versus time accelerometer response. The power of the sampled signal as derived from the power spectral density, in the frequency bands of particular interest (0-100 Hz, 100-300 Hz and 500-700 Hz) and the total power of the sampled signal. The power spectral density is calculated by Welch's method for average periodograms of overlapped, windowed signal sections, based on a discrete-time Fourier transform of the samples of the process using a Fast Fourier Transform (FFT) algorithm [9].

The surface vibration measures which are specifically related to surface vibrational events are as follows. Mean and standard deviation of the amplitude and mill rotation phase of surface vibrational events; Mean and standard deviation of the phase angle weighted by the amplitude of surface vibrational events; Histograms of surface vibrational event magnitude. Total number of events and the ratio of large to small

amplitude events. Contour plots of vibrational event numbers as a function of SAG mill phase angle and event amplitude. Event amplitude and phase angle (at the maximum positive excursion of the accelerometer response for the event).

All the above measures are derived for the total and each revolution of a continuous monitoring period. The signal analysis software has been implemented in the MATLAB programming environment.

RESULTS

A typical accelerometer response trace as a function of time for a rotating SAG mill is shown below (Fig. 1). In this case four full mill revolutions of data were recorded. There is clear evidence of periodicity in the amplitude of events registered by a single accelerometer as a function of mill rotation angle.

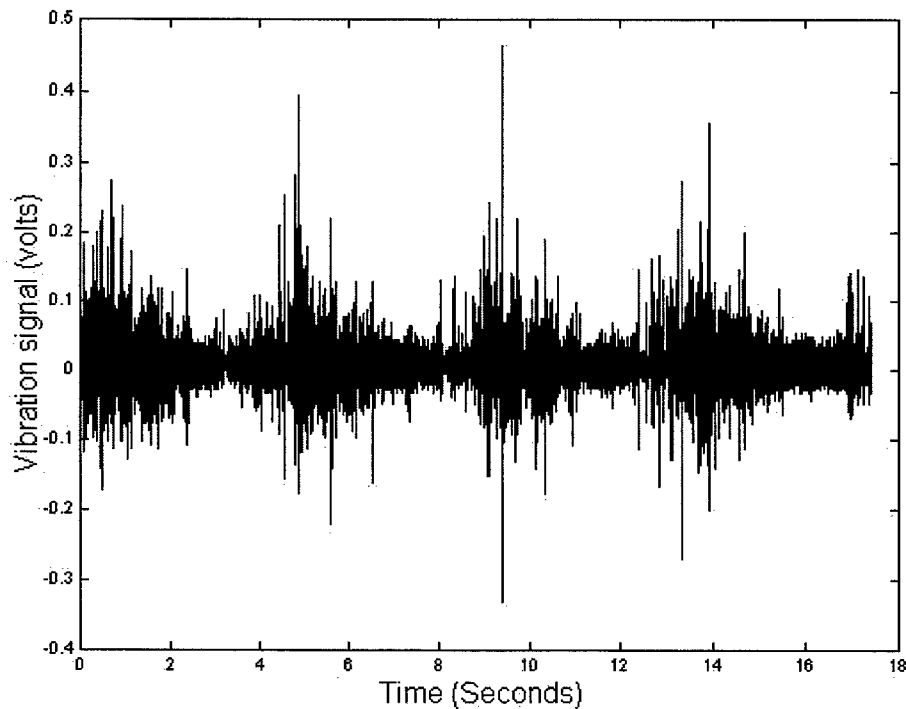


Fig. 1. Accelerometer response on a rotating SAG mill as a function of time.

A contour plot of vibrational event numbers as a function of SAG mill rotation phase angle and event magnitude allows vibration events to be identified with particular locations in the rotation cycle of the mill. In this manner differences may be identified in SAG mill operation both between rotation periods and with changes in mill operating conditions. The position of the average or event amplitude weighted average phase angle of acoustic events is a quantitative measure in this regard.

Figure 2 shows a vibrational event polar contour plot for an 'average' mill revolution, based on the same data set as the Fig. 1 accelerometer response. There is clear evidence of the expected localisation of large events in regions where the charge is thought to be in contact with the liner. Hence it is inferred that there is strong damping of vibrational waves as they propagate around the shell. Conversely, Fig. 2 also shows that there is a registration of lower strength vibrational events in regions where the charge and grinding media are not expected to be in contact with the shell liner. These lower amplitude signals are most likely due to surface vibrational waves propagating around the shell from other regions. However, it seems likely that very high energy events registered by an accelerometer on the outside of the mill shell at a particular phase angle do reflect collisions in the adjacent region of the inner liner. The identification of the boundaries of contact of the SAG mill charge with the liner is potentially important as they may be used to deduce the volume of the charge and investigate lifter bar design efficiency. The ability to discriminate large energy

impacts above the 'toe' of the charge (position where cataracting media impact on the charge) could be used to infer the liner wear rate.

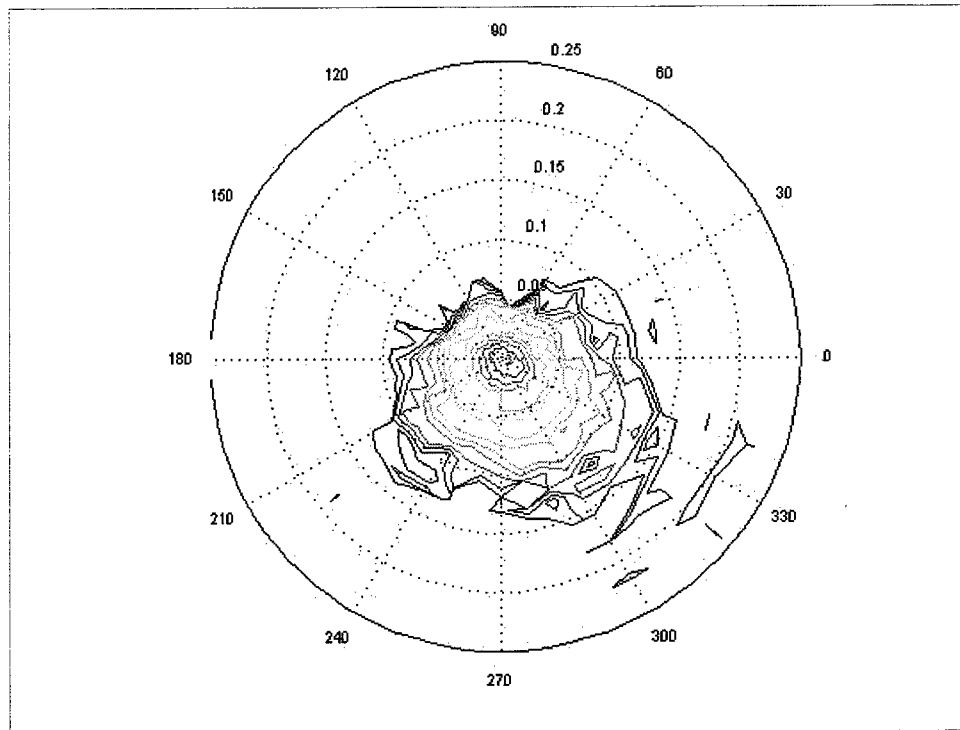


Fig. 2. Polar contour plot of the natural logarithm of the number of vibrational events ($e^{0.5}$ number intervals) as a function of amplitude (volts) in the radial direction and mill rotation phase angle (degrees anti-clockwise from the 3 o'clock position) in the azimuthal direction. Contours of large numbers of events are at low amplitudes. Note that the mill is rotating clockwise.

Figure 3 shows the amplitude weighted average acoustic event phase angle over a revolution as a function of revolution number. As expected from examining Fig. 2, the average event phase angle is in the quadrant associated with the 'toe' of the charge. There are clear differences between the average phase angle for the different operating conditions previously mentioned. The average position angle associated with low feed rate conditions is significantly less than the same measure for high feed rates. This is physically reasonable as one would expect the steady state volumetric loading to be less under low feed rate conditions. The average phase angle for a high rotation rate is clearly larger than the corresponding measure for a low rotation rate. Again, this is physically reasonable as one would expect a high rotation rate to result in more grinding media impacting higher up the liner wall at a higher rotation angle. Significant changes in average event phase angle occur with mill rotation number. This may be indicative of bulk 'sloshing' of the charge on frequencies equal to or lower than the mill rotation frequency.

Figure 4 shows PSD plots for extremes in the mill rotation and feed rate operation variables. Spectral features that are sensitive to mill operating conditions are apparent for frequencies < 100 Hz, around 100-200 Hz and near 600 Hz. Surface vibration power is higher at low tonnages for all frequencies shown. High mill speeds result in increased surface vibration power at low frequencies (<~100 Hz). The total power associated with the signal PSD is significantly higher for both the low tonnage (~40%) and high speed (~30%) operating conditions in comparison with values for the respective extremes in operating conditions. Sharper spikes in the power spectrum can be seen at frequencies below ~500 Hz in the case of high mill rotation speed. All these features are well above the background noise level and are postulated to be related to features of the charge motion.

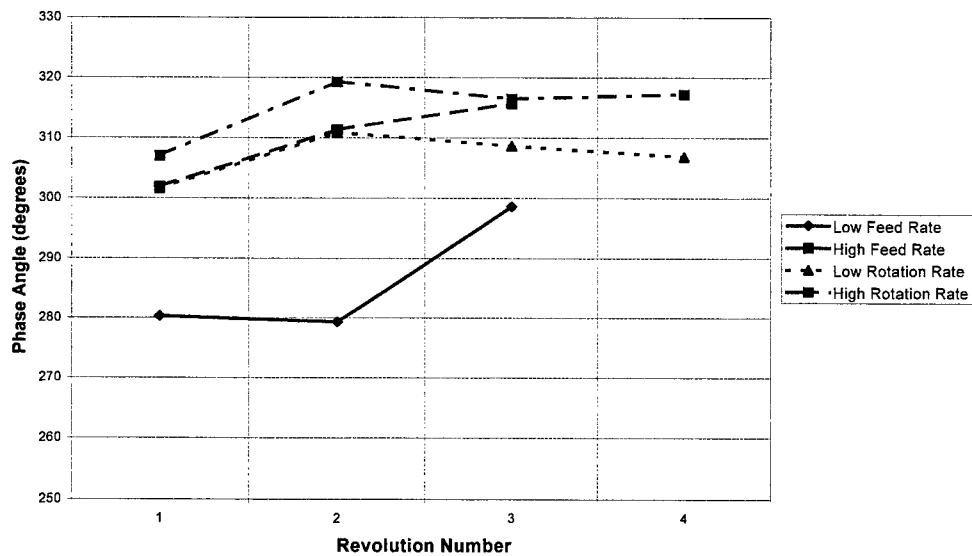


Fig. 3. Amplitude weighted average event phase angle as a function of revolution number. Low and high feed and rotation rates.

In Fig. 4 the frequency range is restricted to a maximum of 1000 Hz, based on experience that ~80% of the surface vibration signal power is within this range. However, PSD plots have also been obtained up to the Nyquist frequency (half the sampling rate) and show a prominent spectral feature at a relatively high frequency (~18000 Hz) that is also sensitive to mill operating conditions (particularly high rotation speed).

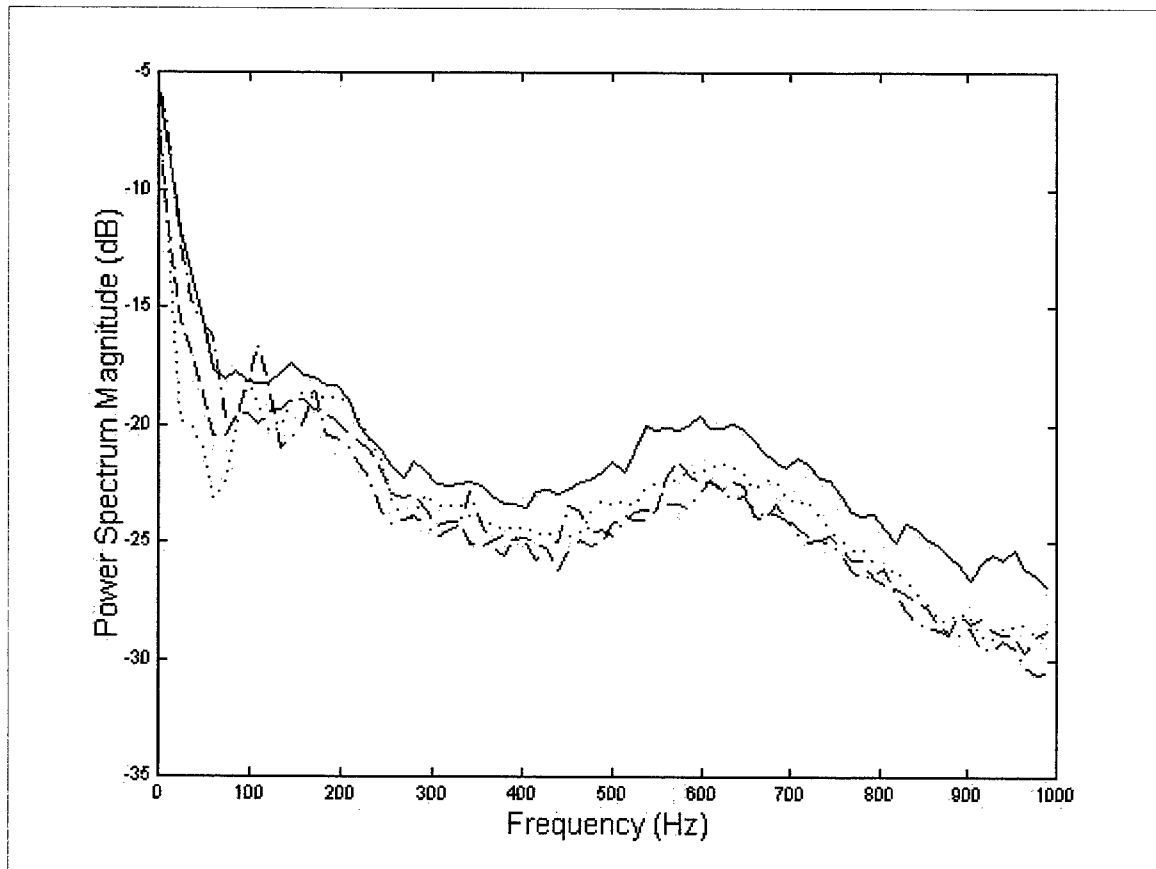


Fig. 4. PSD plot up to 1000 Hz (2^{12} FFT length and Hanning windowing, 2^{11} number of samples overlap, and nil detrending - see [9]). low tonnage: solid line, high tonnage: large-dashed line, low speed: dotted line, high speed: dot/dashed line.

Histograms of the number of vibrational events as a function of the amplitude of the accelerometer response associated with the event are another way of demonstrating both the total number and distribution with amplitude of vibrational events as a function of mill operating conditions. In this manner operating conditions that lead to a relatively large number of very high amplitude vibrational events can be easily identified. Such cases may correspond to conditions of high liner wear. Conversely, concentration of events at a relatively low amplitude may indicate ineffective particle grinding within the SAG mill. For the extremes in operating variables mentioned, the number of acoustic events per unit time significantly increased with tonnage (~2425 - 2550 per second) and speed (~2350 - 3350 per second). However, the associated average amplitude of events decreased with increasing tonnage. Hence high tonnage conditions may imply more events per unit time but with less average energy per event. High speed conditions imply more events per unit time but in the conditions investigated, slightly less average energy per event.

It was thought that the standard deviation of the sampled signal would be another useful measure of activity in the mill at a given set of conditions. Figure 5 shows a plot of sampled signal standard deviation and mill gross power for three different mill speed settings.

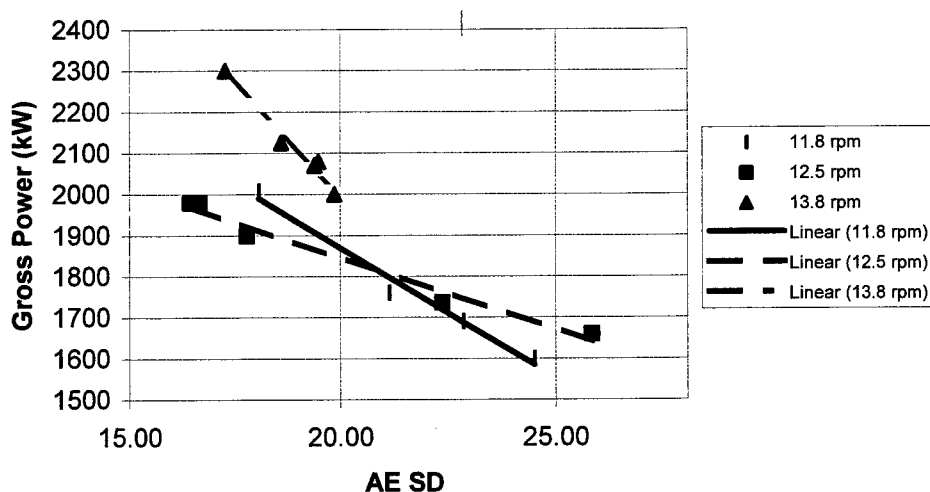


Fig. 5. SAG mill gross power as a function of the standard deviation of surface vibrations for three different mill speed settings.

The mill gross power corresponds to the overall load level in the mill (and feed rate at steady state conditions) i.e., the higher the gross power, the higher the load level. Other manipulated variables, namely ball addition and pulp density, were held constant. At each speed, the relationship between standard deviation and gross power is reasonably fitted by a linear function with negative slope. An increase in the mill gross power (load level) leads to a decrease in the standard deviation of the signal. This result is consistent with experience where the mechanism at play is thought to be increased damping as a result of the higher load level in the mill. The fitted equation for each speed differs and is likely to be indicative of a combination of other effects. Similar relationships are being developed for other surface vibration measures previously listed. The relationships shown in Fig. 5 mean that for a given speed, the measured standard deviation relates directly to gross power of the mill which itself is an indirect measure of the load level. The sampled signal standard deviation increases with mill rotation speed. Evidence has also been found for a decrease in surface vibration standard deviation with an increase in pulp density and an increase of the same with grinding ball addition. These results are to be expected in terms of the ability of the pulp to dampen collisions between the grinding media and the liner and the numbers of grinding media impacts.

CONCLUSION

The operation of a SAG mill results in the generation of high frequency surface waves on the outside of the rotating shell due to collision events involving the charge of ore particles, grinding media and sometimes the mill liner. Acoustic emission (surface vibration) monitoring, using an accelerometer attached to the outside of the rotating shell, has been performed on an industrial SAG mill. Vibration data were collected for a wide range of mill operating conditions. Results for dynamic steady states were analysed by a variety of signal processing techniques and a number of characteristics of the signals were found to be sensitive to changes in mill operating conditions in a physically reasonable manner.

Surface vibration analysis shows considerable promise as a non-invasive means of SAG mill process and condition monitoring. The technique may eventually be used as an input to control systems, an indicator of liner wear rate and lifter bar efficiency. Further work is proceeding on very low frequency monitoring, multiple accelerometer arrays to investigate signal localisation and coupling of results to SAG mill grinding models. Potential applications of the surface vibration technique extend well beyond SAG mills to other mineral processing equipment and in fact any machine that processes material and has a requirement for a better understanding of the mechanisms both from a processing and condition monitoring viewpoint.

ACKNOWLEDGEMENTS

The authors wish to thank CSIRO for permission to publish this paper.

REFERENCES

1. Kelly, E.G., Spottiswood, D.J., 1982. Introduction to Mineral Processing. John Wiley & Sons.
2. Wills, B.A., 1992. Mineral Processing Technology. Pergamon Press.
3. Jaspan, R.K., et al., 1986. ROM mill power control using multiple microphones to determine mill load, Proceedings of the Gold 100 Conference, SAIMM, Johannesburg.
4. Watson, J.L., Morrison, S.D., 1986. Estimation of pulp viscosity and grinding mill performance by means of mill noise measurements, Minerals and Metallurgical Processing, 3, Nov., 216.
5. Watson, J.L., 1985. An Analysis of Mill Grinding Noise, Powder Technology, 41, 1, 83-89.
6. Watson, J.L. and Morrison, S.D., 1985. Indications of Grinding Mill Operation by Mill Noise Parameters, Particulate Science and Technology, 3, 1.
7. Hou, R., Williams, R.A. and Hunt, A., 1998. Acoustic Monitoring of Hydrocyclone Performance, Minerals Engineering '98 Abstracts, 53-55.
8. Boashash, B., Hornsby, D.T. Iskander, D.R., 1998. On-Line Monitoring of Dense Medium Cyclones in Coal Preparation Plants Using Vibration Signal Analysis, XIII International Coal Preparation Congress Proceedings, Australian Coal Preparation Society, II, 469-478.
9. Kay, S.M., 1988. Modern Spectral Estimation. Englewood Cliffs, NJ:Prentice Hall.

Novel Polymeric Electrochemical/Chemical Sensors and Display Devices Integrated with Artificial Intelligence

A. Talaie***, J.Y.Lee***, Y.K. Lee****, J. Jang*, D.J. Choo*****,
S.H. Park*****, G. Huh*****, J.A. Romagnoli**

* Physics Department, Kyung Hee University, Dongdaemoon-ku, Seoul, 130-701, Korea

** Chemical Engineering Department, Sydney University, Sydney NSW 2006, Australia

*** Chemistry Department, NSW University, Sydney, Australia

**** Information Display Department, Kyung Hee University,
Dongdaemoon-ku, Seoul, 130-701, Korea

***** Chemistry Department, Kyung Hee University,
Dongdaemoon-ku, Seoul, 130-701, Korea

Email (corresponding author): aus100@hotmail.com

ABSTRACT

In this study we report on the use of novel and intelligent polymers in various applications in the areas of intelligent sensing and smart devices. We address the reproducibility and irreversibility issues in sensing technology by integrating with a computer and using micro-devices. Examples of a pH sensor, electronic tongue and a display device are demonstrated.

INTRODUCTION

The scope of chemical/electrochemical sensing and display technology has improved markedly in recent years, particularly with the advent of chemically-modified electrodes, polymeric display devices and electrochemical sensors [1-3]. However, their poor selectivity, repeatability and/or reusability have hindered the practical utilization of these sensors/devices. This is due in part to the tendency for all solid surfaces, including polymers, to undergo irreversible changes that can affect selectivity and reusability. Also of concern is a lack of compatibility between the dynamic nature of these surfaces and the usually adopted passive analytical approach.

Owing to the instability of the sensor response quantification can often only be accomplished by use of either calibration curves or standard addition approach. While these approaches may yield useful quantitative data, they do not fully utilize the capability of the sensors. More significantly, the adoption of these quantitative approaches defeat some of the purposes of modern sensing technology in terms of speed, repeatability, reusability and ease of use. These problems require the adoption of novel strategy for the fabrication of sensors with artificial intelligence that will enable the identification, characterization and classification of the response pattern. Such pattern recognition approach, if feasible, will enable reliable determination of the concentrations of analytes and, thus, reduce the emphasis or concern on the variation of sensor response with repeated use.

The principles of "artificial intelligence" methods based on statistical pattern recognition are similar to those used in human decision making. This usually involves collection of information for a known set of cases (e.g. known concentrations or analytes) in an analogous manner to the experience gained from a given process. However, instead of adopting a subjective approach in synthesizing the information it is replaced by equations derived from the data and used to classify cases into groups. Hence, unlike human decision making, a pattern recognition approach is objective and reproducible. The development of appropriate artificial intelligence methods that can use comprehensive information from the sensor response rather than just the signal magnitude will therefore provide a considerable progress towards realization of the optimum performance of these new devices.

ELECTROCHEMICAL/CHEMICAL SENSORS

Chemical/electrochemical sensors are small-sized devices capable of continuously and reversibly reporting chemical/electrochemical information. In fact, chemical sensors are devices that transform chemical information into analytically-useful signals. The information being transformed ranges from the concentration of a specific sample component to the pH of the operational environment [4].

In a chemical sensor, the main functions required of an intelligent dynamic polymer system are to be able to sense appropriate stimuli and respond to them in a predetermined manner. This is critical to development of the sensor. So, there are two main elements which play key roles in the sensing process of a polymer-based chemical sensor. These are the sensing element and the responsive element (the actuator). The chemical nature of the polymer-operational environment interface (sensing element) is the main factor in controlling the nature of output signals from the responsive element. The chemical nature can be adjusted by using the appropriate polymer. The morphology and other physical properties of the polymer film can change the direction of the chemical/physical information which flows through the polymer. In this study we report on the fabrication of a variety of polymeric materials for different sensing applications.

pH sensor

In the pH sensing area we found a particular polymeric composite structure (Polyaniline-polypyrrole composite) as the best sensing material (Fig. 1).

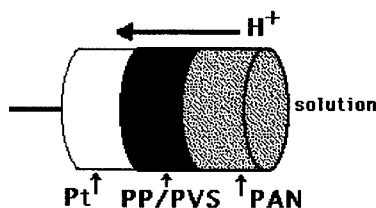


Fig. 1. Conducting polymer composite electrode: Pt (Platinum), PP/PVS (polypyrrole polyvinyl sulphonic acid), PAN (polyaniline)

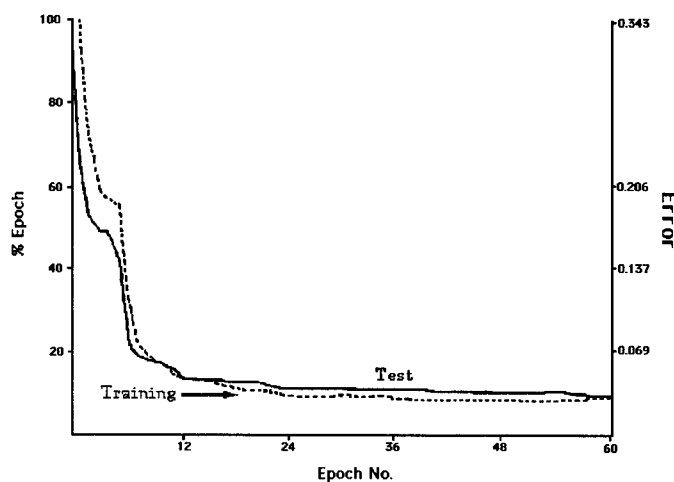


Fig. 2. Test and training profiles during an ANN modeling

Due to lack of reversibility in the pH results, we used several computer algorithms to control the sensor performance. We found the ANN (artificial neural network) method [5] as the optimum computer algorithm amongst others [6] with the best match between prediction and experimental results. As can be seen from Fig. 2, using an ANN algorithm, the training and test processes during modeling meet each other at the same satisfactory level of error. This means that the model represents the system very well.

Electronic Tongue

The most sophisticated sensing systems are found in the human body. For example we can taste via living polymer interfaces. The cellular processes in our tongue are regulated by cell walls comprising dynamic

macromolecules that are able to sense specific chemical stimuli. Amongst these stimuli, chemical ones are of interest to us. There are two types of chemical sensory systems that use different receptors and process information at different locations in the brain: Taste, in which the receptors are specialized sensory cells, and smell, or olfaction, where the receptors are neurons.

Taste receptors of fish are the most sensitive chemoreceptors known. The taste receptors, or taste buds, are not located in the mouth, as in humans, but rather are scattered over the surface of the fish's body. These taste buds are exquisitely sensitive to amino acids. A catfish, for example, can distinguish between two amino acids at a concentration of less than 100 micrograms per liter. The ability to taste the surrounding water in this way is very important to fish to enable them to sense the presence of food in an often murky environment. One group of chemoreceptors in human beings is also concerned with special taste buds, but is not as sensitive as in fish. The taste buds in human beings are located in the mouth (Fig 3). Each taste bud is associated with an afferent neuron. Humans have four kinds of taste buds, each of which respond to a broad range of chemicals. The stimuli to which different kinds of taste buds respond are salty, sweet, sour, and bitter. Our complex reception of taste is composed of different combinations of impulses from these four types. The chemoreceptors are concentrated on different parts of the tongue, sweet and salty at the front, sour on the sides, and bitter at the back (Fig 3a.).

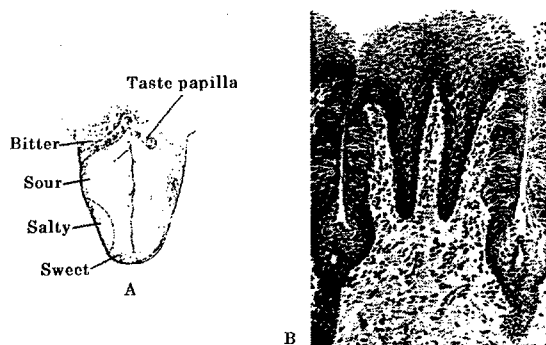


Fig. 3. a. Four kinds of taste buds located on different regions. **b.** Individual taste buds are bulb-shaped collections of chemicals receptor cells that open out into the mouth through a pore.

Learning from the human tongue, our approach is to develop a dynamic sensing electrode by which we can detect solutions which contain salt such as sodium chloride (a salty-taste solution). The salty taste is related to type of salt and its extent depends on the concentration of the salt in the operational solutions. This approach can be used in the food industry to control the extent of salty taste in products if an array of microelectrode is combined with a computer for on-line data processing and prediction (Fig. 4). The typical example of an off-line salt detection (CaCl_2 , LiCl and NaCl) is reported in Table 1.

Table 1. Data for modeling and pattern recognition experiments in off-line salt detection
Due to the nature of the software, the salt type was introduced to the computer as combinations of 0 and 1).

Input Data								
....	156	157	158	159	160	Ca	Li	Na
....	-0.09	-0.09	-0.09	-0.09	-0.10	1	0	0
....	-0.11	-0.09	-0.10	-0.10	-0.11	0	1	0
....	-0.10	-0.10	-0.10	-0.11	-0.10	0	0	1
Output Data								
Salt	Output 1		Output 2		Output 3			
Na	-0.01768		-0.11035		0.95254			
Ca	1.011793		0.073661		-0.10175			
Li	-0.00987		0.98108		0.00729			

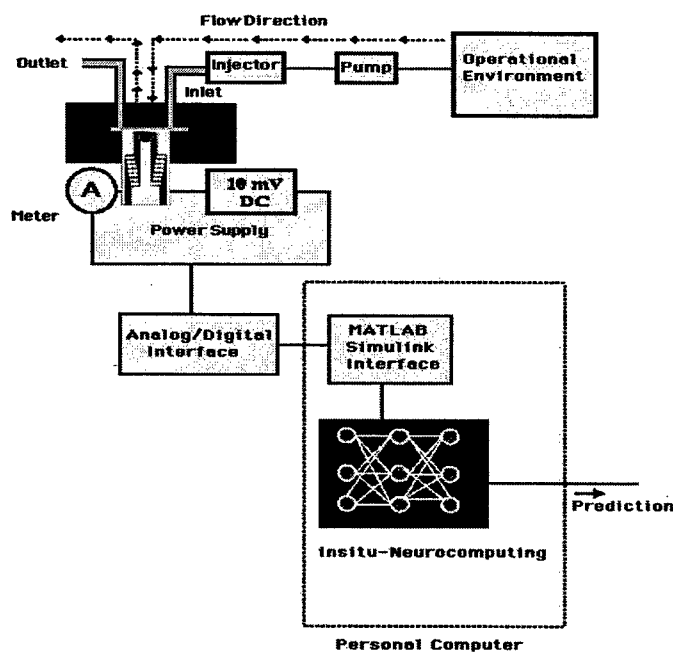


Fig. 4. Diagram of an intelligent device with ability to taste different samples. An array of microelectrodes is connected to a powerful computer equipped with advanced algorithms.

DISPLAY DEVICE

Conjugated polymer electroluminescence devices (CPELD) consist of organic thin layers that are essentially insulating materials. The operating mechanisms involve injection of electrons and holes from cathode and anode to the organic emitter layer, and hole/electron recombination that generates light emission. A typical configuration of a display device (Fig. 5) and the yellowish red color emitted from the polymer (Fig. 6) with its photoluminescence (PL) intensity are reported below. The next step in this research is to integrate the computer with powerful algorithms to this device to control the intensity and the emitted color using different polymeric layers.

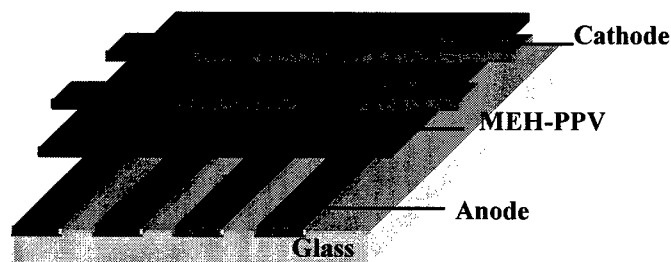


Fig. 5. Configuration of polymeric display device using MEH-PPV (Poly(2-Methoxy,5-(2'-Ethyl-hexyloxy)-P-Phenylenevinylene)

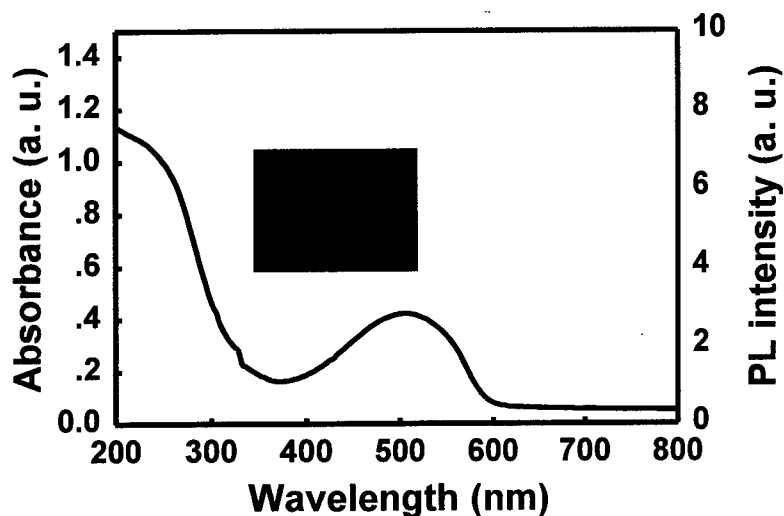


Fig. 6. Typical absorbance and PL intensity profiles for a polymeric display device material

CONCLUSION

The new systems presented in this paper show examples of integration of data and models for chemical/electrochemical sensing and display applications. These integrated systems enable us to control the dynamic nature of devices. Although the system suggests promising industrial applications it needs further development in its neuro-computing step for an *in-situ/real time* modeling and prediction.

REFERENCES

1. A. Talaie, 1997. Conducting polymer based pH detector: A new outlook to pH sensing technology, *Polymer*, 38, 1145-1150.
2. Y. Ikariyama, W.R. Heineman, 1986. Polypyrrole as a detector for electroinactive anions by flow injection analysis, *Analytical Chemistry*, 58, 1803-1808.
3. A. Talaie, Y.K. Lee, J. Jang, D.J. Choo, S.M. Park, S.H. Park, G. Huh, I. Lee, 1999. The effect of transition metals on the PL intensity of MEH-PPV, *Thin Solid Film*, submitted.
4. A. Talaie, 1996. Conducting polypyrrole and polyaniline-based chemical sensors, *Chemistry in Australia*, December, 570-572.
5. A. Talaie, A. J. Romagnoli, 1996. Application of artificial neural networks to the real time operation of conducting polymer sensors: A pattern recognition approach, *Synthetic Metals*, 82, 27-33.

MATERIAL PROPERTIES UNDER DRAWING AND EXTRUSION WITH CYCLIC TORSION

L.X. Kong, P.D. Hodgson, L. Lin and B. Wang

School of Engineering and Technology, Deakin University,
Geelong, Vic 3217, Australia

ABSTRACT

Many cold and hot worked metals undergo strain softening and hardening when subjected to cyclic plastic deformation. The degree of strain softening depends on the amount of prior cold work or heat treatment and upon the strain magnitude and cycles applied. In this work, a lead bar was extruded and a copper bar was drawn through a cyclically twisting die in a specifically designed experimental rig. The drawing/extrusion load fluctuated at the same frequency as that of die twisting. The maximum load was equivalent to monotonic deformation when the die was changing direction. The degree of the reduction in load for both the drawing and extrusion depends on the deformation conditions and requires optimisation for the process.

INTRODUCTION

Controlled cyclic plastic strain experiments have been conducted by many investigators generally to determine the resistance of metals to fatigue failure. In the course of those experiments, the resistance to deformation is measured during the cyclic loading as a function of the number of stress cycles. From this, curves can be constructed of the stress amplitude, or the stress range, versus the number of cycles, the cumulative plastic strain, or some other appropriate measure of the strain accumulation. It has been found that for materials which are in the soft or annealed state and are capable of significant strain hardening in ordinary tension tests, cyclic hardening occurs as a result of the controlled cyclic plastic strain. However, this cyclic hardening is quite different from monotonic hardening where the deformation occurs in one direction only. Under monotonic straining, hardening increases monotonically with strain such that a relationship can be formulated using many constitutive models such as unified dislocation density based constitutive model developed by Kocks, Estrin and Mecking [1-3].

When the material is initially in a heavily cold worked state prior to application of cyclic plastic strain, it has been found [4] that with subsequent controlled cyclic plastic strain, cyclic softening occurs. That is, with each further application of strain the stress reached at the point of reversal is less than that of the previous reversal. For either the annealed or prior cold work initial states, a stabilised hysteresis loop is eventually established which suggests that regardless of the initial condition of the material, the stress range reached after cycling by controlled plastic strain is the same [5]. Based on this experimental observation, Kong et al [6, 7] developed a constitutive model to model the cyclic strain hardening and softening.

Due to the significant change in the properties brought about by the application of cyclic plastic strain following severe cold working, a process which can induce mechanical annealing will have great technological significance. Among possible implementations of the technique into industrial applications, Korbel and Bochniak [8] developed an approach called the Structure-Based Design of Metal Forming Operation (SBDMFO), aiming to develop a metal forming process to make the plastic forming easier and cheaper. Cywinski and Wusatowski [9] conducted an experiment to examine the effects of strain path change using tension-torsion-tension tests and found that the application of the torsion at the late stage of the draw or tension test increases the softening induced by cyclic deformation. This coincides with Coffin's observation that strain softening depends on the degree of the cold working prior to the cyclic straining[5].

In the current work, an experimental rig was constructed capable of conducting both extrusion and drawing with minor modification to study the straining hardening and softening by introducing cyclic deformation.

The rig is similar to a real process to assess the potential for industrial implementation. The experimental results for both extrusion and drawing are analysed to evaluate the straining softening under different deformation conditions.

EXPERIMENTAL

The schematic of the experimental rig for both drawing and extrusion is shown in Figure 1. There are three stages with the drawing process. Die 1 and 3 are fixed. The second die can rotate along its axis. The material is hardened in the first die. It then experiences a cyclic straining applied by the second die. The specimen at the second die is therefore subjected to a combined tensile and cyclic torsional strain. The specimen may be work (cyclically) hardened or softened depending on the reduction applied in the first stage of the drawing. At the third stage, the specimen is subjected to further cyclic straining because the material is cyclically twisted before entering the die.

The extrusion is the same as traditional extrusion except that the die can be rotated in both directions. After the specimen fills the container by plastic deformation during the initial compression, the friction between the specimen and the container prevents the upper section of the specimen within the container from rotating with the die while the lower section twists with the die as it is extruded through the die. The material approaching the die is therefore subjected to a combined deformation of extrusion and cyclic torsion. Unlike the drawing process, there is only one die for the extrusion process. During the extrusion, the material undergoes a complex deformation which is not as easily classified as that of drawing. When the material is within the container, it is subjected to a small amount of monotonic extrusion and is slightly work hardened. When the material is in the die, the material is initially monotonic and cyclic hardened and then possibly cyclic softened depending on the deformation achieved.

The transmission mechanism for both the drawing and extrusion is the same. A control system was used to reverse the turning direction and adjust the twisting frequency and angle. As the experiment was set to investigate the effect of cyclic straining, the drawing/extrusion velocity remained the same for all the test. The load was recorded against the displacement with a strip chart recorder and a PC data acquisition system at a frequency of 10 Hz.

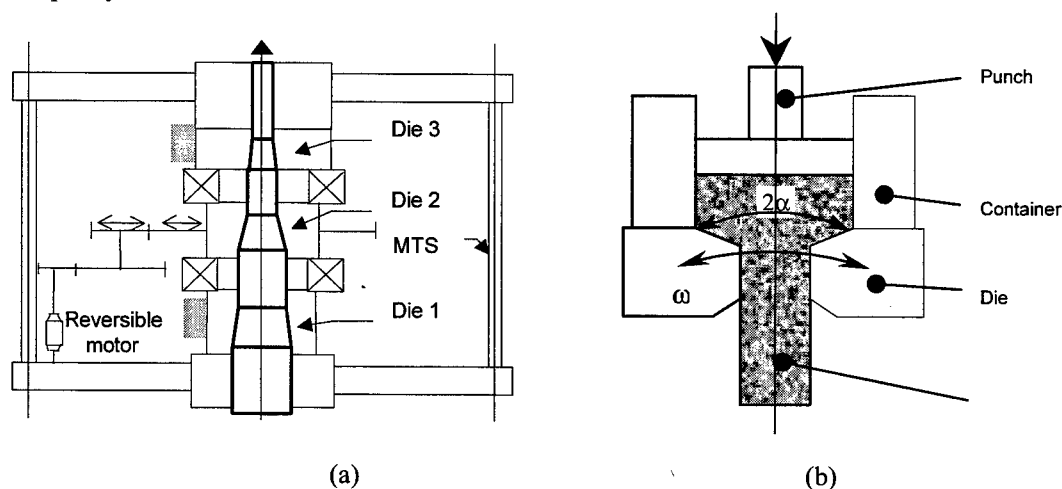


Fig. 1. Experimental scheme: (a) drawing [6] (b) Extrusion [7]

EXPERIMENTAL RESULTS AND DISCUSSIONS

Extrusion test

The material used in the extrusion test was a commercially pure lead. The lead bar was cast with a diameter of 40 mm and extruded to the specimens with a diameter of 16 mm and a length of 65 mm. The specimen was then annealed at 160° C for 2 hours to produce a uniform grain size and consistent properties. The tests were carried out on a screw-driven tension-compression machine with torsion applied by twisting as

discussed above. Experiments were performed at ambient temperature using the conditions listed in Table I. The diameter of the specimens was reduced from 16 to 8 mm corresponding to an area reduction of 75%.

Table I: Experimental scheme of the extrusion test

	Die semi-angle	Die rotation speed	Twisting frequency
1	90°, 60°, 45°	15°	20 cycles/min
2	60°	7.2°, 15°, 30°	15 cycles/min
3	60°	30°	20, 10, 6.7 cycles/min

Monotonic extrusion was firstly performed. The extrusion force varies with punch travel in a typical way with three clearly distinguished phases [10]. In the first phase, the workpiece expands and fills the container as the punch starts compressing the workpiece through the die. The extrusion force initially increases rapidly in this phase, achieves a recognisable peak value and then drops dramatically. In the steady-state phase, the extrusion force fluctuates and, generally, decreases gradually due to reduction in friction between the workpiece and container. However, the decrease in extrusion force is very smooth. Once the workpiece advances to the unsteady-state phase, the extrusion force decreases in a more rapid rate and then increases dramatically as the workpiece is completely extruded through the die and the punch starts to contact the die.

If the extrusion process is accompanied by cyclic torsion, the workpiece is then subjected to a combined deformation which leads to a different pattern in the extrusion force and punch travel relationship. Fig. 2 superimposes the effect of cyclic torsion on the extrusion force for a lead bar twisting at a rate of 10 cycles/min. Although the three regimes observed and used for monotonic extrusion can still be seen, the peak extrusion force in the first regime is not as easy to extract as in the monotonic extrusion. This may suggest that a steady-state phase can be more easily achieved with the introduction of cyclic torsion. However, of all the effects, the most pronounced is reduction in extrusion force during twisting of the die.

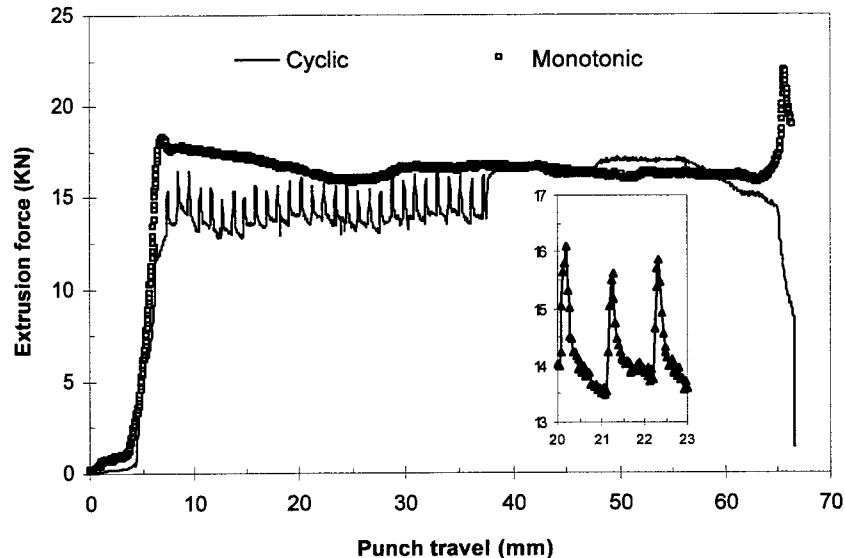


Fig. 2. Extrusion force versus punch travel for extrusion with cyclic torsion at frequency of 10 cycles/min

In each cycle, variations in the load during extrusion with cyclic torsion was divided into three stages regardless of the deformation conditions: softening stage, steady stage and hardening stage (Fig. 2)[7]. The steady stage can still be a softening stage but softening is very slow compared to the first stage. If the average or cyclic extrusion force is set as the boundary of steady state, the transition points of these three stages are easily defined (Fig. 2).

Using the same experimental scheme, the effect of die semi-angle, die rotation speed and die twisting frequency were studied. It was found that a die with a semi-angle of 60° will achieve the maximum cyclic softening, resulting in the lowest extrusion force [11]. Therefore, the die with a semi-angle of 60° was considered as the optimal die and was used for further experimentation. To eliminate the difference introduced by the experimental conditions, the tests were performed on the same specimen; for example by changing the twisting speed from $7.2^\circ/\text{s}$ to $15^\circ/\text{s}$ and then to $30^\circ/\text{s}$ during extrusion when evaluating the influence of die rotation speed. Figure 3 shows the extrusion force during extrusion through the die. It is observed that the maximum extrusion force during the three phases remains almost the same and therefore changing the twisting speed does not reduce the extrusion force at the turning points where the motor, and hence die, reverse direction. This confirms our previous observation that, at the moment that the die stops during reversing direction, the extrusion temporarily becomes monotonic deformation. In contrast, the degree of cyclic softening increases with the twisting speed. As shown in Figure 3(b), the extrusion force, at the steady state of every cycle, reduces uniformly with an increase in the die twisting angle. When the twisting speed is $30^\circ/\text{s}$, the cyclic softening is approximately 20% compared to a monotonic deformation; at twisting speed of $7.5^\circ/\text{s}$, the cyclic softening is only about 6%.

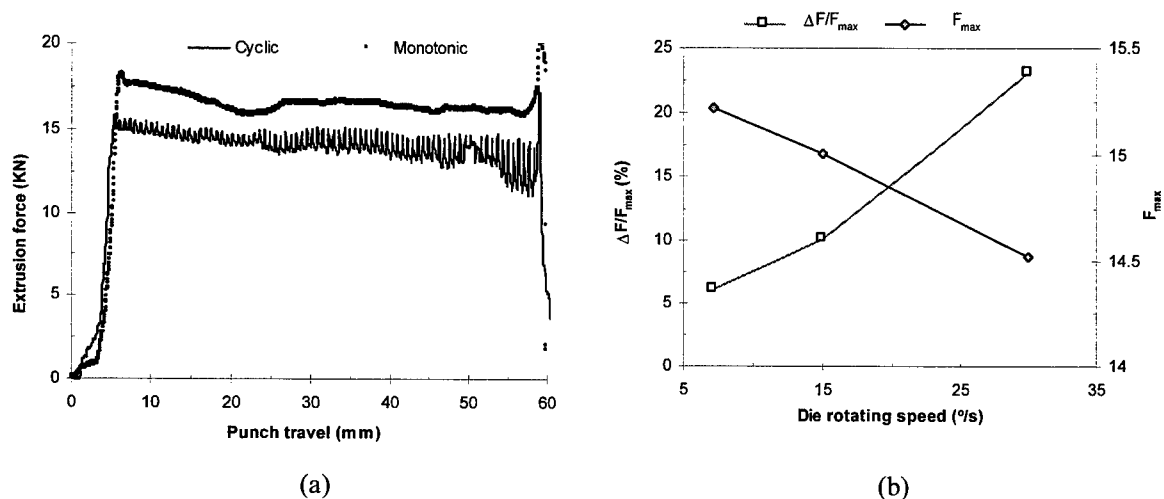


Fig. 3. Extrusion with different die rotation speeds: (a) experiment (b) analysis

From Fig 3, it seems that a faster die twist rate leads to a higher degree of cyclic softening. However, experiments involving changing the twisting frequency do not support this. As the die twisting range remains the same for all the tests, a higher twisting frequency leads to a faster die rotation.

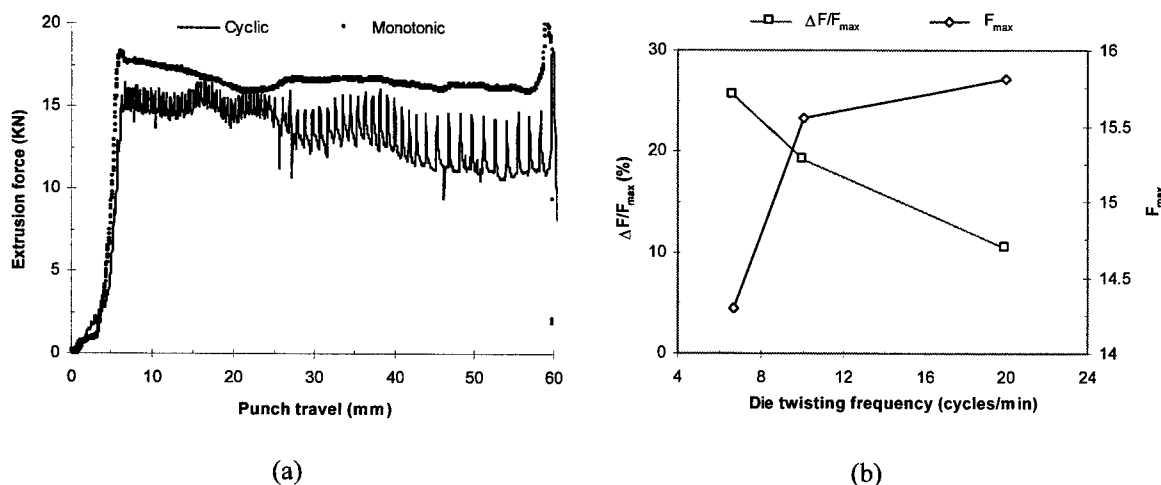


Fig. 4. Extrusion with different die twisting frequencies:
(a) experiment (b) analysis

When the die twists at a frequency of 20 cycles/min, the degree of cyclic softening is only 10% (Fig. 4). As the twisting frequency decreases, the degree of cyclic softening increases dramatically. When die twisting frequency reduces from 20 cycles/min to 10 cycles/min, the degree of cyclic softening increases from 10% to 20%. Further reduction in twisting frequency to 6.7 cycles/min leads to further softening up to 25%. The experimental result further verifies the previous observation that the steady state in every cycle is enlarged if the die twists at a slower speed because the time to achieve steady state depends on the extrusion speed rather than the other deformation conditions. In addition, the change in die twisting frequency results in little variation in the maximum extrusion force (Fig 4).

Drawing test

A commercial copper was used for the drawing tests. The experimental scheme is listed in Table II. The diameter of a copper bar reduced from 9.6mm to 8.6mm with an area reduction of 20%.

Table II: Experimental scheme of the drawing test

	Drawing speed (mm/min)	Die rotation speed
1	100, 50, 20	30°
2	20	20°, 30°, 40°, 60°, 80°

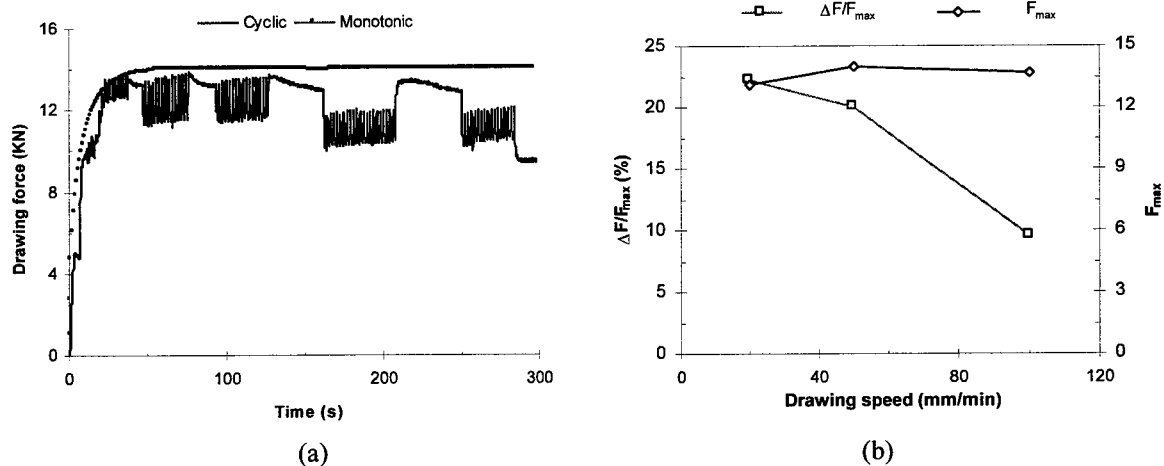


Fig. 5. Drawing with different drawing speed: (a) experiment (b) analysis

The drawing force against drawing time is shown in Fig. 5 where only the drawing speed varied. A similar phenomenon to that observed in extrusion tests can also be seen in drawing tests. When the die starts twisting, the drawing force reduces from the monotonic drawing force. The degree of the reduction or softening caused by twisting die depends on the experimental conditions, ie drawing speed here. The experiment is repeatable as the test on drawing speeds of 50 and 20 mm/min was carried twice following a short period of monotonic drawing. If the drawing speed is the same, the minimum drawing force reduced to the same level. The maximum drawing force and the degree of the cyclic softening was analysed using the same technique for extrusion. The maximum drawing force is more stable than the extrusion tests, as there are no different stages observed in extrusion force. However, the softening is more obvious, particularly if the drawing speed is at a low level (ie 20mm/min). When the drawing is at 100mm/min, the degree of cyclic softening is only about 10%. As the drawing speed reduces, the degree of the cyclic softening increases to about 22% for the drawing speed of 20mm/min.

Although the drawing speed has a large influence to the degree of cyclic softening, varying die rotation speed does not affect both the maximum drawing force and the degree of the cyclic softening significantly (Fig. 6). The only exception is the rotation speed of 20°/s where the minimum drawing force is higher than others, which means a lower level of softening is achieved. The stability of the drawing force suggests that

the variation in the die rotation speed does not lead to a remarkable change in the drawing force and the selection of this parameter depends only on the experimental or working conditions most applicable.

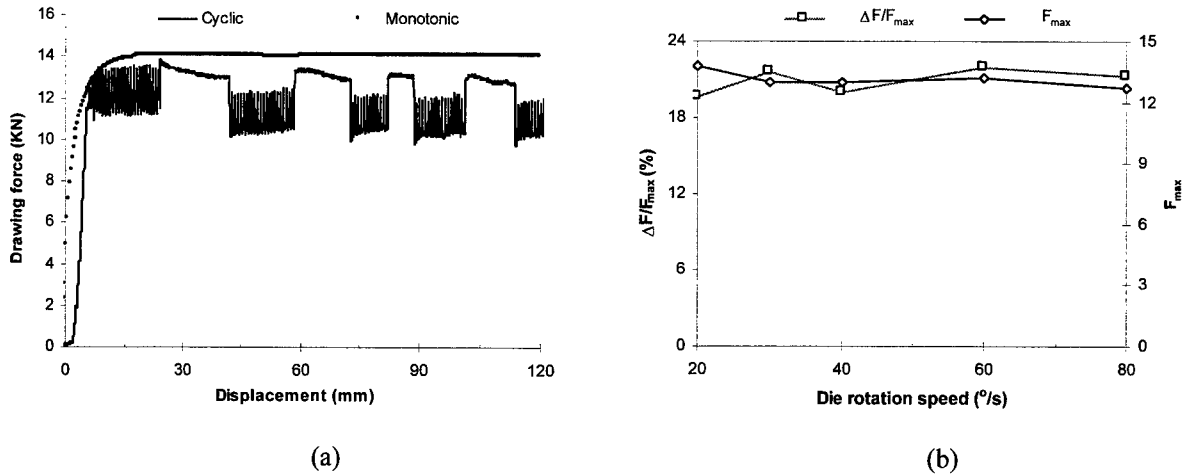


Fig. 6 Extrusion with different die rotation speed: (a) experiment (b) analysis

CONCLUSION

Lead was tested under combined extrusion and cyclic torsion while copper was drawn with a die twisting. Different experimental conditions were used to examine the process. Regardless of the extrusion or drawing process, the application of a cyclic torsion leads to a cyclic softening. However, the softening occurs only when the die is turning. If the die stops twisting or when the die is changing direction, the extrusion/drawing force increases to the monotonic force. It is observed that variation in experimental conditions generally lead to a change in extrusion/drawing force and in the degree of cyclic softening. A clear understanding of the process with cyclic strain path change requires further experimental and analytical study to more materials and deformation conditions.

REFERENCES:

1. Y. Estrin and H. Mecking. 1984. A unified phenomenological description of work hardening and creep based on one-parameter models. *Acta Metall.*, **32**, 57-70.
2. F. Kocks and H. Mecking. 1981. Kinetics of flow and strain-hardening. *Acta Metall.*, **29**, 1865-1875.
3. L. X. Kong and P. D. Hodgson. 1999. Improving the Prediction Accuracy of Constitutive Model with ANN Models. IPMM'99, Hawaii.
4. D. S. Dugdale. 1959. Stress-strain cycles of large amplitude. *J. Mech. and Phys. Solids*, **7**, 135.
5. L. F. J. Coffin. 1967. Cyclic strain-softening effects in metals. *Transactions of the ASM*, **60**, 160-175.
6. L. X. Kong, P. D. Hodgson, and B. Wang. 1999. Development of constitutive models for metal forming with cyclic strain softening. *Journal of Materials Processing Technology*.
7. L. X. Kong. 1999. Constitutive Modelling of Extrusion of Lead with Cyclic Torsion. *ASME Journal of Engineering Materials and Technology*.
8. A. Korbel and W. Bochniak. 1995. The Structure Based Design of Metal Forming Operations. *Journal of Materials Processing Technology*, **53**, 229-237.
9. Cywinski M and Wusatowski R. 1994. The Influence of Controlled Twisting on Drawing of OFHC Copper and Armco Bars. *Journal of Materials Processing Technology*, **45**, 299-304.
10. K. Lange. 1985. *Handwork of Metal Forming*. Dearborn, Michigan: Society of Manu. Engineers.
11. L. Lin, L. X. Kong, R. Bathgate, P. D. Hodgson, B. Wang, and G. Lu. 1999. The Effect of Die Design on the Bar Drawing with Combined Torsion. Tooling'99, Melbourne, Australia.

Artificial Neural Networks II

Neural Network Method for Inverse Modeling of Material Deformation

Nenad Ivezic, John D. Allen Jr., and Thomas Zacharia

Computer Science and Mathematics Division
Oak Ridge National Laboratory
Oak Ridge, TN 37831, USA

ABSTRACT

A method is described for inverse modeling of material deformation in applications of importance to the sheet metal forming industry. The method was developed in order to assess the feasibility of utilizing empirical data in the early stages of the design process as an alternative to conventional prototyping methods. Because properly prepared and employed artificial neural networks (ANN) were known to be able to codify and generalize large bodies of empirical data, they were the natural choice for this application. The product of the work described here is a desktop ANN system that can produce in one pass an accurate die design for a user-specified part shape.

INTRODUCTION

Sheet metal forming applications have traditionally involved costly, iterative, and time-consuming design methodologies based on physical prototyping. Recently, finite element analysis techniques have been used with varying success to model and predict sheet metal deformation processes during the design stages [1]. Additionally, certain methods from the field of artificial intelligence (AI) have been brought to bear on the problem of automating design of tooling for manufacturing metal stampings [2]. Yet ANN methods, among the most powerful AI tools, have barely been explored in this context, although they have been used to solve engineering design problems involving identification, learning, and prediction of scalar or vector design quantities [3, 4, 5]. In contrast, the method described in this paper attempts to capture deformation relationships governing 3-D mappings in sheet metal stamping applications. Presented results demonstrate that the inverse of these relationships can be captured and used to predict die geometry that will, under specified material and stamping conditions, allow production of a user-specified part shape.

DEVELOPMENT OF THE METHOD

Realization of an ANN-based die design tool depended critically on successful solution to several inter-related problems. Central to these was development of a suitable method to represent die and part shapes to the input nodes of a neural network (or sometimes, multiple networks). For network training, these data would be of two classes. Input data would represent the conformation of the part for which a die geometry was sought while output data (those data applied to the output node(s)) would specify to the network, the result that should be obtained for the corresponding inputs. These data sets would represent in some manner, the die geometry that the network must learn to produce in response to the input data. Although, in principle at least, a network should be able to learn in the context of almost any representation, relatively simple and often unobvious variations in the representation scheme can produce striking differences in the efficiency with which learning takes place. Described below are several of the surface-representation methods investigated during development of the die prediction system.

Simplest in concept, but surely least effective of all, was a direct mapping of part and die data onto the input and output nodes of a very large network. Here, the input and output nodes were arranged in "rectangular" fashion, the number of each node class being equal to the size of the part/die meshes (never more than 25 by 25 for the first attempt). Although this method was adequate for initial demonstration of

the general technique, it was doomed to failure by the impossibility of obtaining enough training data to condition the network properly. Even if sufficient data were available, it is doubtful that any single network architecture would be able to generalize the shape transformation relationships over a useful spectrum of part shapes under such a representational scheme. If it could do so at all, network prediction performance would be very sensitive to minor displacements of presented shapes relative to the input grid.

A method based on two-dimensional Fourier transforms derived its utility from the fact that a relatively small fraction of the total number of Fourier components (of the order of 10 %) could represent a 3-D shape with considerable fidelity. As employed to represent die and part geometries to a neural network, the transform was applied to each of the data sets to be presented for evaluation. The training goal for the network was the development of the functional relationships linking the Fourier components representing a part configuration (together with the associated data defining material and forming parameters) with the corresponding Fourier components representing the die that would, under the specified set of material and forming parameters, be suitable to produce that part. When a fully-trained Fourier-based network is presented with a data set representing a part (and the associated material and forming parameters) for which a die design is required, the network produces at its output nodes a set of values representing the scaled Fourier components of the desired die. Recovery of the die specification only requires an appropriate renormalization and retransformation into a Cartesian coordinate representation.

A somewhat similar method based on two-dimensional Wavelet transforms was developed as the next step in the search for an efficient representation method for surface data. The 2D Wavelet method, like the Fourier method described above, derives its utility from the fact that a very small fraction of the total number of wavelet components (perhaps 1-2 %) can represent a surface with considerable fidelity. For a reasonably narrowly defined set of shapes, it is observed that the identity of the important Wavelet components is fairly consistent across all members of the set with variations among set members being represented principally in amplitude differences among the important components. Substantive variations among set members will be represented by the appearance of dominant wavelet components more or less singular to the set members bearing those variations. It is the combination of consistent and singular components that is captured and represented to the network for subsequent processing.

Although 2-D Wavelet representation reduced the number of data required for shape representation by a greater factor than that obtaining for the 2-D Fourier method, the reduction still appears insufficient to support completely general network training. The single greatest advantage is one of speed. Once trained, the network executes in milliseconds on any reasonably modern desktop computer.

The most successful, investigated methods for die/part representation is the so-called Weighted Patch method. This scheme was predicated on the assumption that shape gradients near a point on a material surface are more likely to be predictive of the response of material at that point to deformation forces (imposed, for instance, by a punch and die) than are shape gradients at positions more removed from the point. Equally reasonable is the assertion that "near" effects must be represented more completely to the neural network than more "distant" ones if the model is to capture important metal forming relationships.

In the Weighted Patch Method, a fixed pattern of "averaging regions" is scanned over the data array representing a part for which a die is to be developed. In a typical implementation, there may be 19 averaging regions of graduated sizes (9 for each of two orthogonal axes and one central element, the shortest being one array element square, the longest of the order of 8 to 10 array elements in length by one element in width) comprising a cross shaped "Patch Geometry" of the general form suggested by the schematic representation of Figure 1. These averaging regions are distributed along orthogonal axes (axes typically parallel the length and breadth of the part). Each averaging region is represented at the network input layer by a single node.

The training goal for a Patch Method network is to produce, for each element of the part array, the value of the difference between part elevation and corresponding die elevation at that point. When the network

is fully trained a predicted die configuration can be trivially derived from these differences for any presented part shape.

The Patch Method confers the advantage that a somewhat more accurate representation of the part surface is retained in the reduced data employed by the network than is the case with either the Fourier or Wavelet data representation methods. Perhaps more important than this is the fact that, since each point on a surface is treated as a separate training example, the network can develop a much better codification of relationships linking part and die shape than can a network employed for either of the previously discussed methods.

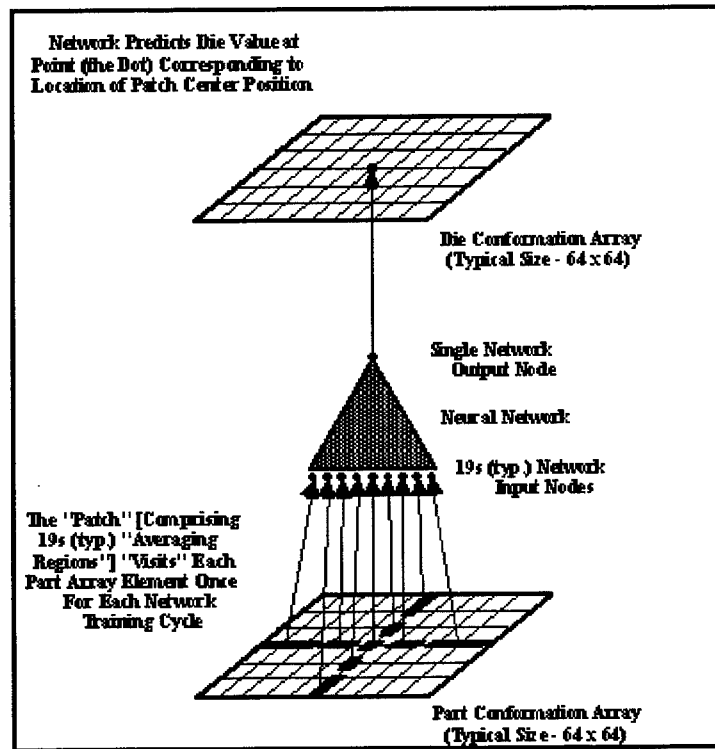


Fig. 1. Patch method schematic diagram (parameter nodes are not shown)

The penalties paid for realizing these advantages are two in number. First is the requirement that the network "visits" each of the elements of the part definition array many times during training. For the 64×64 data sets employed to develop the results presented here, there are 4096 such elements, each of which may be visited several hundred times during a typical training session. The second penalty is the memory intensive character of the method. For each of the 4096 potential patch-center locations, it is necessary to store the 19 values comprising the Patch (for a total of 19×4096 , or 77824 values) for each data set to be processed by the network (of which data sets there may be hundreds).

EXPERIMENTAL DEMONSTRATION

The potential power of the Weighted Patch method is well-illustrated by the results of a simple experimental demonstration, for which the goal was to demonstrate that a network trained under the method could generalize its training results well enough to produce a die design for a member of a shape class to which it had not been exposed during training.

For this demonstration, part definitions were created from corresponding die definitions by a relatively simple deformation model that included two variable parameters, one representing a material property, the other a forming variable. The characteristics of the model were chosen to ensure that the resulting shape transformations would be of sufficient complexity to challenge the network.

Model output was derived from die representations of the three generic shapes depicted in Figure 2. For each die shape class, variation of the relevant defining parameters resulted in a spectrum of die geometries from which corresponding parts would be "formed" by the deformation model. Pan dies, for instance, were defined by mean base width, mean base height, mean top width, mean top height, top and bottom aspect ratios (so that either base, top, or both could assume trapezoidal shape), center elevation, and top tilt angle. Additional parameters specified the presence or absence of ridges and their spatial extents relative to the base die dimensions. Similar parameters (with the addition of eccentricity and appropriate omissions) defined the model Top Hat and Ellipsoidal Shell dies. The value of one last variable determined the extent to which die corners were rounded. In all cases, the die elevations were defined on a regular, 64 by 64 element, square grid. Although this representation may at first appear contrary to the conventions of standard deformation modeling, it is entirely in keeping with the metrological considerations that would dictate the collection of data from physical die/part pairs.

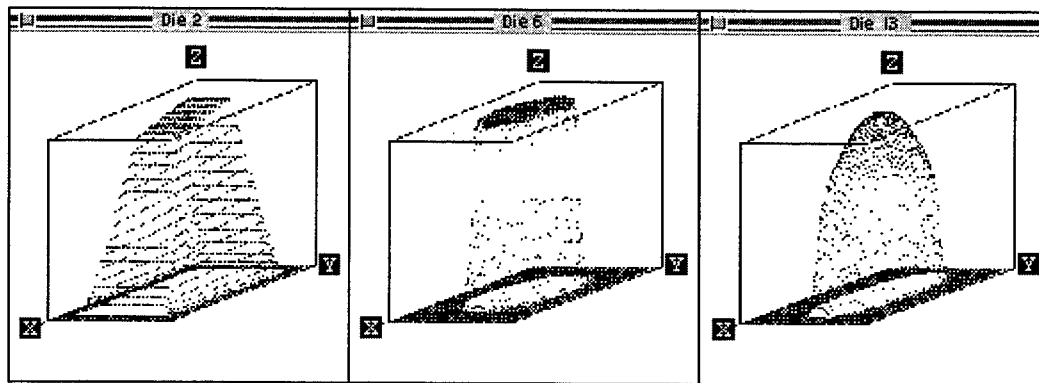


Fig. 2. Pan, TopHat, and Shell Die Shapes

The simple deformation model created part data from each of the die representations as a parameterized displacement function of the local and the global geometries of the die. It should be noted that although the simple deformation model calculates material displacements in each of the three axes, the final part data are interpolated back onto a regular 64 by 64 grid. These, expressed as an array of elevations, are stored with the generating die specification (again, an array of elevations) and a list of the "material" and "forming" parameters associated with the part/die pair to form a single training data set.

Network training involved 40 part/die pairs. Of these, 20 defined variations on the Pan geometry. The remaining 20 defined variations on the Elliptical Shell geometry. For the Pan parts, elevation, aspect ratio, top tilt, and top and bottom shape (square or trapezoid) were varied. Elevation values (normalized to mean base size) fell in the range (0.062 - 0.125). The mean top size to mean base size ratio varied between 0.5 and 1.0. Aspect ratio (mean length to mean width) varied over the range (0.47 - 1.0). Top tilt angle values (in radians) fell in the range (0.0 - 0.5). Where they were relevant, the same parameter ranges characterized the generated Elliptical Shell part/die pairs. Parameters of the deformation model were established to limit elevation differences between the die and predicted part to less than ~10 %.

Pre-processing of data included the interpolation and normalization steps required to produce the 4096 input Patch vectors and the corresponding output value for each of the 40 training data sets and the single TopHat test data set. Node components and output values for the training data were normalized to the range (-0.85 , +0.85). Post-processing to recover the die shape involved an inversion of the output normalization process and subtraction of the resulting values from the corresponding stored part elevations (since the network is trained to reproduce point-wise differences between part and die).

The neural network architecture employed in this experiment was a conventional Perceptron of four layers. The input layer comprised 21 non-bias nodes, the first hidden layer 13 non-bias nodes, the second hidden layer 3 non-bias nodes, and the output layer 1 non-bias node. Of the 21 input nodes, 19 were employed for representing Patch data, two for representing the values of the "material parameter" and "forming parameter" values stored with each data set. The network was trained according to the gradient-descent, backpropagation of error method. During early stages of training, the common practice was to process multiple networks simultaneously. When a clearly superior network emerged, the poorer members of the initial set of networks were excised. During final training, learning rate and momentum were slowly relaxed from initial values of .1 and .001 respectively to .001 and .00001.

Figure 3a is the conventional scatter diagram and represents network training and testing results as a function of expected result. It is emphasized that network predictions for the single test case (4096 more points) are included in the scatter diagram (and, surprisingly, account for very little of the scatter). At first glance, it might appear that the trained behavior of the network is insufficiently accurate to be of much utility for die shape prediction. Figure 3b, a three-dimensional perspective representation of the network-predicted die, suggests otherwise.

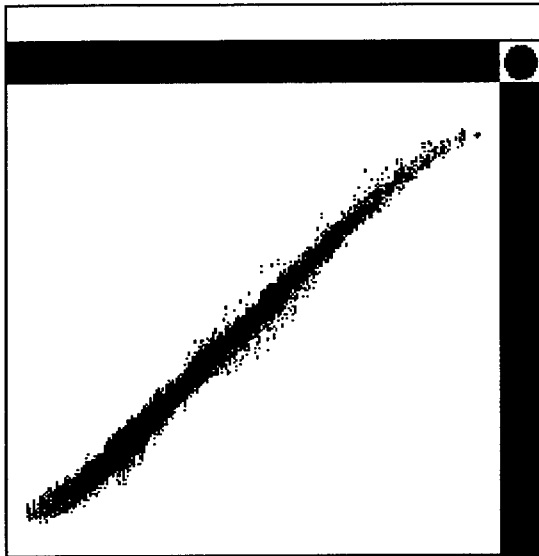


Fig. 3. a. Scatter diagram for a training process. Note: this plot includes 163840 points for the Pan and Shell training data and 4096 points for the single TopHat test case.

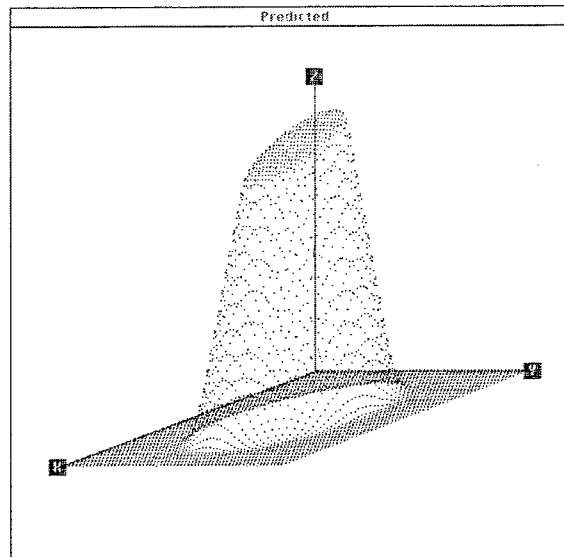


Fig. 3. b. 3-D perspective plot of predicted die.

Figure 4 provides a composite view of the Patch Method results and depicts the original part (dashes), the true die (i.e., the die from which the part was generated by the deformation model) (solid), the network-predicted die (dots) and the prediction error (times 10) for each point in the cross-section (crosses). In all cases, the prediction error has proved to be no greater than one part in 256 of the presented part elevation.

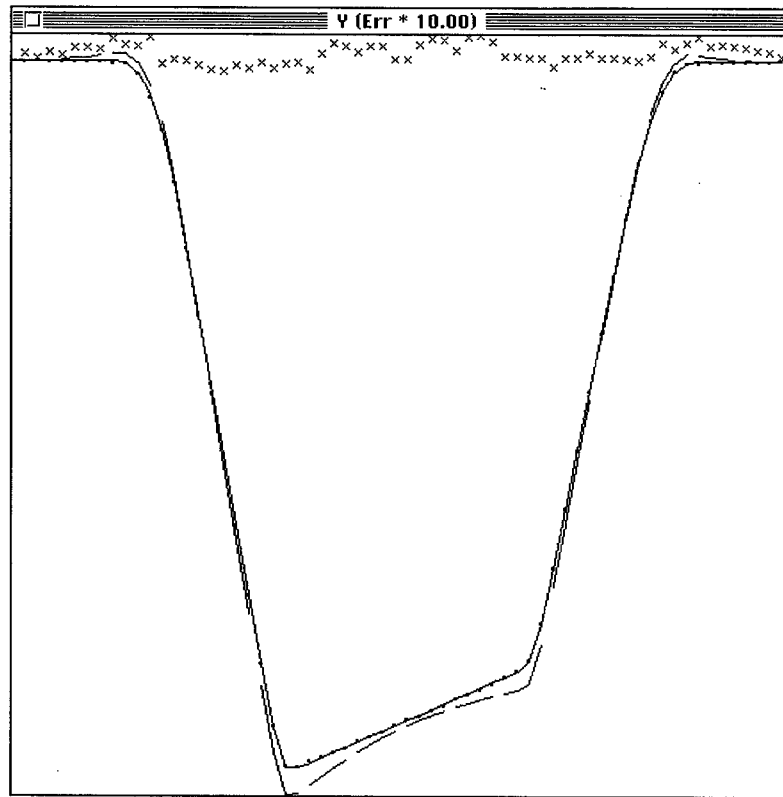


Fig. 4. Cross-section of part, (dashes), predicted die (dots), true die (solid), and die/true-die error x10 (crosses).

SUMMARY

In this paper, development and experimental results were described of a neural network-based system able to inverse-model material deformation. Several alternate representation methods were discussed. Of these, the Patch Method proved to be the most accurate to predict die geometries from user-supplied part shapes. It was used to obtain experimental results to demonstrate the potential utility of the method.

ACKNOWLEDGEMENTS

This research was sponsored by the Laboratory Technology Research Program, Oak Ridge National Laboratory, managed by the Lockheed Martin Energy Research Corporation for the U.S. Department of Energy, under contract number DE-AC05-96OR22464. This submitted manuscript has been authored by a contractor of the U.S. Government under Contract No. DE-AC05-96OR22464. Accordingly, the U.S. Government retains a non-exclusive, royalty-free license to publish or reproduce the published form of this contribution, or allow others to do so, for U. S. Government purposes.

REFERENCES

1. S. Yang, K. Nezu, 1998. Application of an inverse FE approach in the concurrent design of sheet stamping. *Journal of Materials Processing Technology*, 79, 86-98.
2. B.T. Cheok, A.Y.C. Nee., 1998. Trends and developments in the automation of design and manufacture of tools for metal stampings. *Journal of Materials Processing Technology*, 75, 240-252.
3. N. Ivezic, J.H. Garrett, Jr., R. Ganeshan, 1991. Generalized Hopfield Network for Structural Optimization. in *Proc. Artificial Neural Networks in Eng. (ANNIE91)*, St. Louis, Missouri, 849-854.
4. N. Ivezic, J.H. Garrett, Jr., 1994. A Neural Network-based Machine Learning Approach for Supporting Synthesis. *Artificial Intelligence for Eng. Design, Analysis, and Manuf.*, 8, 143-161.
5. N. Ivezic, J.H. Garrett, Jr., 1998. Machine learning for simulation-based support of early collaborative design. *Artificial Intelligence for Eng. Design, Analysis, and Manuf.*, 12, 123-139.

An Adaptive Artificial Neural Network to Model a Cu/Pb/Zn Flotation Circuit

Saiedeh Forouzi ^{**}; John A. Meech ^{*}

^{*} Department of Mining and Mineral Process Engineering,
University of British Columbia, Vancouver, BC, Canada

⁺ currently with Minnovex Technologies Inc., Toronto, Ontario

Email: forouzis@minnovex.com, jam@mining.ubc.ca

ABSTRACT

In this paper, we describe the planning and development of an Artificial Neural Network model of line 3 of the Copper/Lead flotation circuit at Brunswick Mining's concentrator at Bathurst, New Brunswick. The prototype model predicts the copper and lead assays of the concentrate streams of this rougher flotation circuit. In the model, the actual values and rates of change in the main process variables such as head grades, reagent addition, mass flow, density, pH, temperature, cell level and grind size are treated as inputs. The global error in both training and testing of the model is used to indicate the accuracy of the model. The model is fully adaptable, i.e., it can be updated when required to account for ore and/or processing changes that are not currently included in the ANN because of lack of instrumentation or reliability of measurements. The adaptation algorithm is used to select current data to replace records in the existing training and testing datafile. Retraining is conducted whenever the model accuracy declines to a pre-defined target value. The algorithm determines the frequency of retraining. The final system will be expanded to calculate a total of 12 assays using a separate ANN model for each. All models are independently updated. This approach to Artificial Neural Networks provides plant engineers with a process model that is always current and reasonably accurate. Model access provides flexibility in adjusting set-points to achieve increased efficiency in the control of process variables.

BACKGROUND

Brunswick Mine is a large massive-sulfide deposit of Copper (Cu), Zinc (Zn), Lead (Pb) and Silver (Ag), located in northern New Brunswick. Sphalerite, galena, chalcopryrite, tetrahedrite and a large assortment of silver sulpho-salts are the economic minerals. The Brunswick Concentrator produces four concentrates of Zinc (Zn), Lead (Pb), Copper (Cu) and Bulk (Pb/Zn). The mill process consists of the typical unit operations such as crushing, grinding, flotation and dewatering. There are a Jaw Crusher and two Gyratory Crushers underground for preliminary size reduction. The crushing and grinding plants on surface, consist of a Cone Crusher and two Short-Head Cone Crushers, Rod Mills and Ball Mills. In September 1998, an Autogenous Mill was commissioned to replace the crushing plant and rod-mills. This has caused major changes in the process operations. In the Brunswick Mill, grinding is carried out in three parallel lines. Ball milling is done in primary and secondary mills in a closed loop with cyclones. The discharge from Line 1 and 2 are joined together before flotation, whereas the discharge from Line 3 is processed separately. Selective flotation is carried out in two separate stages in Lines 1&2 and Line3. After flotation, all concentrates are dewatered by thickening and filtering and then shipped to smelters. The general process flowsheet after the AG mill installation is shown in Figure 1.

Due to time limitations, the project focused on the Cu-Pb flotation circuit of Line 3. Application of the same procedure to each of the other flotation circuits is the ultimate extension of this project. In the Cu/Pb flotation circuit, the milled product is aerated to promote galena flotation and depress pyrite. Soda ash is added during grinding to adjust pulp pH while xanthate (SIPX & PAX) and AF241 are used as collectors. The circuit consists of a rougher and two cleaner stages. The flowsheet of this circuit is shown in Figure 2. Assay information is available on-line on a cycle time of 5 minutes generated from 5 streams using an on-stream X-Ray Fluorescence Analyser. Plant operators use these XRF assays to adjust reagent additions, water addition, tonnage rates and flotation cell levels to control product assays at pre-defined set points [1].

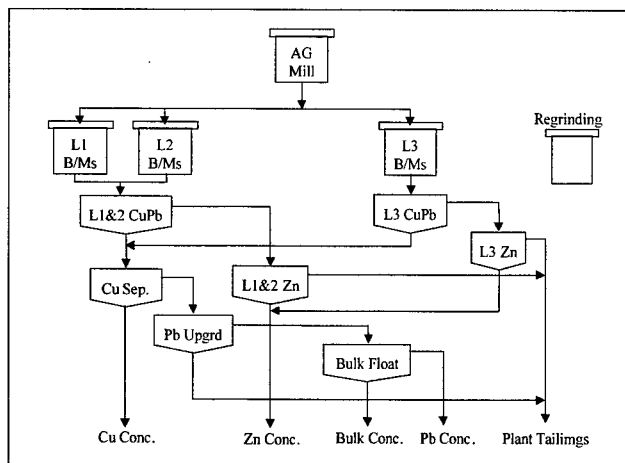


Fig. 1. The general flowsheet of the Brunswick Concentrator after AG mill installation

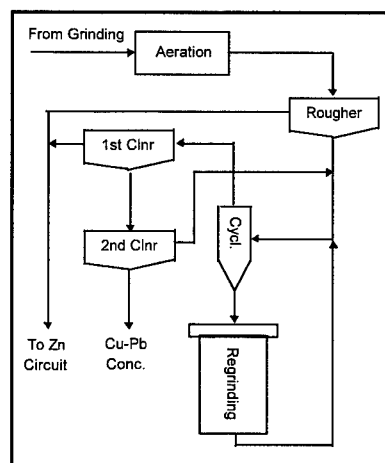


Fig. 2. The flowsheet of the Copper-Lead flotation circuit.

THE PROJECT

Artificial Neural Network (ANN) modeling has advantages over other methods. ANNs need only identify model inputs and outputs. Only general knowledge about the domain is necessary. In ANN modeling, one does not need the form of inter-variable relationships since the algorithm on its own, approximates an appropriate form. An ANN is built by introducing large amounts of historical data to the model to gradually improve its knowledge of process relationships. The main disadvantage is a dependency on the quality of the training data. Providing sufficient valid and accurate data is often difficult.

To implement this project we used software that works with real-time process data, trains and tests an ANN model, infers facts based on rules, executes procedures, and makes decisions in real time. Due to the availability of GenSym's G2, GDA and NeurOn-Line and their past application for process control at Brunswick Mining, we decided to use these tools. G2 provides a foundation to deliver advanced control applications using methods such as neural networks, fuzzy logic, and genetic algorithms. Integrated with G2 and GDA, NeurOn-Line provides a complete development environment to create intelligent real-time applications for on-line process monitoring, optimization, and model-based reasoning. G2 is designed to interface with external programs and plant automation systems such as data bases, PLCs, and DCSs.

Objectives

The Adaptive Artificial Neural Network model built in this project, was designed for future process control. Our main purpose was to have an ANN model that can learn about the relations between process inputs and outputs and how such relationships change over time. On the other hand, a process model enables us to anticipate the effect of any input scenario on process outputs. In this way, a system can be designed to manipulate control variables to compensate for changes in those variables that are beyond our control, such as head grades. The required adjustment is determined by established rules in a knowledge base and from the predicted process outputs from the ANN model. A schematic of this plan is illustrated in Figure 3.

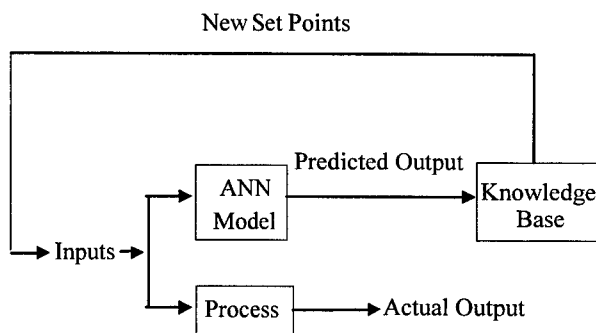


Fig. 3. Schematic Diagram of Compensating Feed-Forward Model-Based Control Algorithm.

ANN ARCHITECTURE

To train the model, the most important variables and their acceptable ranges were identified. The inputs included grades of copper and lead in the Cu/Pb concentrate, the head grades, reagent additions, pulp pH, mass flows, pulp densities and temperatures, grind size, cell levels and rates of change in all the variables. The outputs included change in Cu and Pb grades of the Cu/Pb concentrate. The Cu and Pb grades are applied as inputs to attempt to compensate for variables which are impossible or difficult to measure. The number of nodes used was: 58 in the input layer, 30 in the hidden layer and 1 in the output layer for each of two networks; all nodes were biased. The sigmoid function was applied to reduce the emphasis of outliers. All data were scaled between 0 and 1 using the minimum and maximum values of each variable.

IMPLEMENTATION

Implementation of the project constituted two parts. First, the ANN model was trained using historical data, while the second part consisted of developing an algorithm to decide on the frequency of retraining the model as new real-time data accumulated. The retraining procedures were also designed in this step. The data used to train and retrain have a vital effect on the model accuracy. Before introducing data to the model, the data must be pre-processed. This pre-processing included filtration, synchronization and scaling.

Data Pre-Processing

To avoid noisy data, a suitable range was defined for each input and output. These ranges were determined according to the experience of plant experts and from the historical plant data. Every datum passes through a filter to check its validity in terms of these acceptable ranges. If the datum is valid, it is used in the data set to train the model, otherwise the entire data pair is discarded.

Due to the dynamic nature of the process, there is a delay between a change in any variable and its effect on the output. These phase lags are functions of the circuit flowrate and the effective volumes of pipes, cells and pumps. The phase lag between each variable and output is calculated based on Equation 1:

$$\text{phase lag} = (\text{volume} \times \text{fullness coefficient}) / \text{flowrate} \quad 1.$$

The volume of all pipes and unit operations were determined from drawings, specifications, and direct measurement. The fullness coefficients were defined according to the experience of experts -- as well as measurement of cell levels, etc. being used to modify these delays in real-time. The flow rate at each stage is measured by a sensor in real-time. To apply these delays, the entire circuit is considered as a sequence of individual stages. When a variable is measured, the phase lags between all stages are combined iteratively based on the flow at each stage to synchronize the sensor readings. Due to the variety of ranges and units in the inputs and outputs, the data were scaled between 0 and 1. The scaling method used a linear function so that the maximum and minimum values are 1 and 0 respectively.

Initial Training

For initial training, historical data from July to September 1997 were available for use. This "old" data were pre-processed through filtration, synchronization and scaling. After pre-processing, each data pair was collected into a matrix called a data set in G2. A total of 1300 data pairs were collected. To train the model, the Conjugate Gradients (Fletcher-Reeves) [2] algorithm was used. The data set was randomly divided with 80% used for training and the rest for testing. After several test runs, we decided to train for 150 iterations. In G2, an iteration stands for the number of times the objective function is applied to the entire training data set. Each time the objective function is used, this represents one step. The testing error indicates how well the network fits data not used in training. The fitness option selected was RMSE, i.e., Root Mean Squared Error in which the predicted value is subtracted from the target value to obtain the error.

Retraining

To ensure the neural network model accurately reflects the current process relationships, procedures were developed to check the model against current data and decide if and when the model needs retraining. There are two parameters used to determine this answer: the current root mean square (RMS) error of the model; and the uniqueness of the current data set. (see Figure 4.) For retraining, the updated data set is divided randomly into two separate files for training and testing.

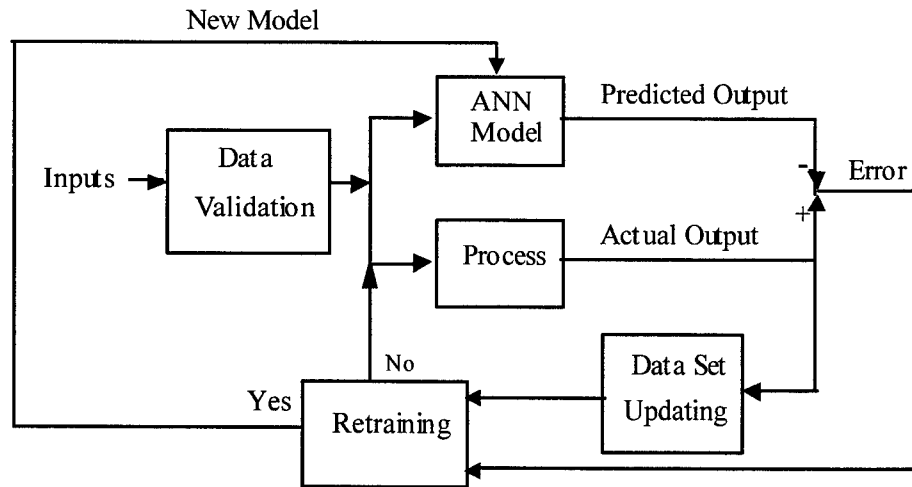


Fig. 4. Schematic diagram of the retraining process.

Parameter 1: The Current RMS Error of the Model

Whenever a data pair is introduced to the trained neural network, a prediction of the process output is made based on the inputs in the data pair. The actual output in the data pair is compared to the model prediction and the relative error is calculated by Equation 2 while the current RMS error is calculated by Equation 3.

$$\text{Relative Error} = (\bar{Y}_k - Y_k) / \gamma_k \quad 2.$$

$$\text{RMS Error of Model} = \left(\left[\sum_{k=1}^n ((\bar{Y}_k - Y_k) / \gamma_k)^2 \right] / n \right)^{0.5} \quad 3.$$

γ_k was set to 100% of each variable's range so that output error would be independent of the actual value.

Parameter 2: The Uniqueness of the Current Data

If a change occurs in a process relationship, the current data pair will reflect this change. By monitoring the current data pair and introducing it to the neural network, those data pairs that represent new relations in the process can be diagnosed. Obviously these particular data pairs should be used to eventually retrain the network. In order to keep the training time relatively constant, we only allow a maximum number of data pairs in the training data set. When a new data pair arrives, an algorithm is applied for data management.

In this algorithm, if the model error is large, the system compares the current data pair with each data pair in the existing data set. If there is an identical one, the current data pair is simply discarded. If the current data pair is unique, it is considered for use in retraining. If the current data pair indicates a changed process relationship, it may be added to the data set or it may replace an old data pair. When the current data pair is totally different from what the model has seen before, it is added to the data set. However if the number of data pairs in the data set is at its maximum level, one of the existing data pairs must be deleted to open space for the new one. When two data pairs are said to be identical or different, the criterion is the relative error. The current data pair with input vector X and output vector Y is compared with row i of the data set with input vector X' and output vector Y' . The relative error for each datum is calculated based on its offset divided by a reference value. The reference values for vectors X and Y are the vectors η and γ respectively.

current data pair	vector X (inputs): $[x_1, x_2, x_3, \dots, x_m]$	vector Y (outputs): $[y_1, y_2, y_3, \dots, y_n]$	
data pair in row i	vector X' (inputs): $[x_1', x_2', x_3', \dots, x_m']$	vector Y' (outputs): $[y_1', y_2', y_3', \dots, y_n']$	
reference value	vector η : $[\eta_1, \eta_2, \eta_3, \dots, \eta_m]$	vector γ : $[\gamma_1, \gamma_2, \gamma_3, \dots, \gamma_n]$	4.

When a current data pair arrives at the algorithm, the Data Comparison Procedure is run. In this procedure, for every j th element of the inputs and every k th element of the outputs, the relative difference and RMS difference between the current data pair and row i of the existing data set is calculated as follows:

relative difference for input j	$e_{xj} = (x_j - x'_j) / \eta_j$	$j = 1, 2, \dots, m$	
relative difference for output k	$e_{yk} = (y_k - y'_k) / \gamma_k$	$k = 1, 2, \dots, n$	
RMS difference for inputs	$E_x = ((\sum (e_{xj})^2) / m)^{0.5}$	$j = 1, 2, \dots, m$	
RMS difference for outputs	$E_y = ((\sum (e_{yk})^2) / n)^{0.5}$	$k = 1, 2, \dots, n$	5.

If any row is found in which the absolute value of all related differences is less than a predefined threshold, then the two data pairs are considered identical and the current data pair is simply discarded:

in row i , $\forall j = 1, 2, \dots, m$ and $\forall k = 1, 2, \dots, n$: if $|e_{xj}| \leq \text{threshold}$ and $|e_{yk}| \leq \text{threshold}$
 then data pairs are identical
 then discard current data pair 6.

Otherwise the current data pair is compared to the next row. This process repeats until all rows are checked. If no row is identical to the current data, the Data Replacement Procedure is called to attempt to find a row where all absolute values of the relative differences in the input vector lie below the threshold while at least one of the absolute values of the relative differences in the output vector lie above the threshold as follows:

in row i , $\forall j = 1, 2, \dots, m$ and $\exists k = 1, 2, \dots, n$: if $|e_{xj}| \leq \text{threshold}$ and $|e_{yk}| > \text{threshold}$
 then current data pair is a new condition
 then current data pair replaces one of the identified rows 7.

If no row is detected, this means there is no data pair in the data set to be replaced with the current data pair so the Data Addition Procedure is called to add the new data pair. If more than one row is detected, then one must be selected to be replaced with the current data pair. The criteria is to select the data pair with the minimum RMS difference in the input vectors and the maximum RMS difference in the output vectors. So to decide between two data pairs with input vector RMS differences of E_{x1} and E_{x2} and output vector RMS differences of E_{y1} and E_{y2} in outputs, at different conditions the selection would be as follows.

- ♦ if $E_{x1} < E_{x2}$ and $E_{y1} > E_{y2}$ then data pair 1 is selected
 - ♦ if $E_{x1} > E_{x2}$ and $E_{y1} < E_{y2}$ then data pair 2 is selected
 - ♦ if $((E_{x1} > E_{x2}$ and $E_{y1} > E_{y2})$ or if $(E_{x1} < E_{x2}$ and $E_{y1} < E_{y2})$ then
 - if $E_{x1} + E_{y1} > E_{x2} + E_{y2}$ then data pair 1 is selected
 - if $E_{x1} + E_{y1} < E_{x2} + E_{y2}$ then data pair 2 is selected
- 8.

The algorithm keeps track of the portion of new data pairs added to the data set since the last training of the neural network model. This indicates the uniqueness of the current data set to that used to train the model.

Accumulation of the Two Parameters

Accumulation of these two parameters determines the need to retrain. Generally when the portion of new data is high, it is expected that the error of the model will also be high or vice versa, but it is possible that the plant may return to a previous relationship and so the correlation can be reversed. Obviously when the error is high and there is a high amount of new data, retraining must occur. Conversely, if the error is low and the data set has not been significantly changed, retraining is not useful. However when both the portion of new data and the RMS model error are located somewhere in the middle of their respective ranges, the decision to retrain is not so easy. The concept is illustrated in Figure 5. To answer the question we used Fuzzy Logic. Membership functions were defined as in Figures 6 a. and b. The decision to retrain is made regarding the variables in each cell of the Fuzzy Associative Map in Table 1.

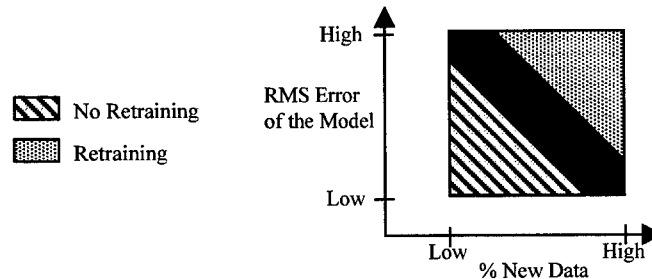


Fig. 5. The decision process to retrain.

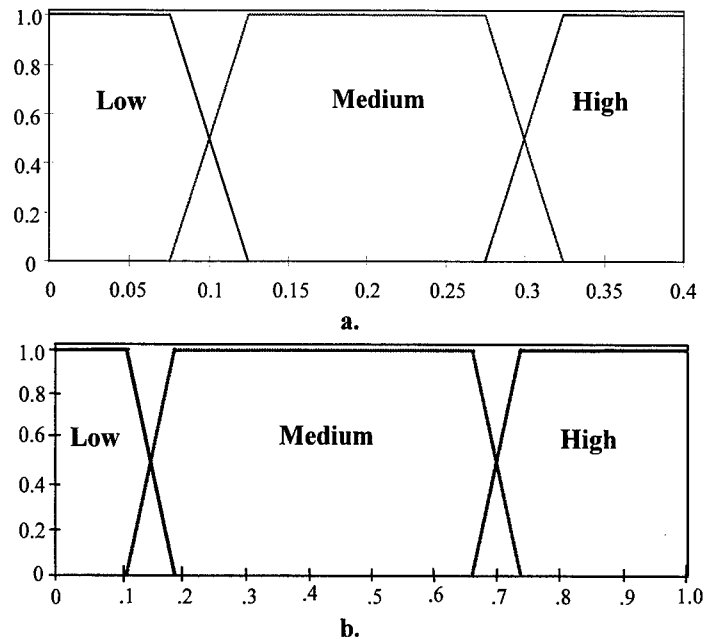


Fig. 6. Membership functions of fuzzy observations for
a. The portion of new data and **b.** the scaled RMS error of the model.

In the case of Potential Training, the requirement for training depends on the slope of the RMS error. Having a positive slope indicates a tendency to move from the fuzzy area of "Potential Training" to the area of "Training Necessary". Conversely, a negative slope shows a tendency to move from the fuzzy area of "Potential Training" to the area of "No Training". Belief in the consequence of all fuzzy rules are accumulated into the fuzzy concepts shown in Figure 7.

Table 1. The fuzzy associative map used to decide on retraining.

Parameters	New Data Portion		
	Low	Medium	High
Low	No Training	No Training	Potential Training
Medium	Potential Training	Potential Training	Training Necessary
High	Training Necessary	Training Necessary	Training Necessary

The need to retrain is resolved by combining the likelihood of the three conclusions of "No Training", "Potential Training" and "Training Necessary". Area-Centroid-Weighting is used to locate the position on the Universe of Discourse and so, the conclusion to be implemented. Belief values below 0.25 are ignored.

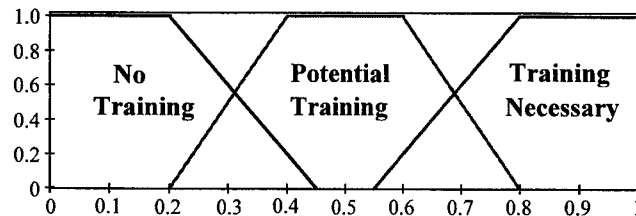


Fig. 7. Membership functions of the consequence of the fuzzy associative map for retraining.

RESULTS AND DISCUSSIONS

When the historical data were collected, several back-propagation neural networks with different architectures were tried. The main parameter varied was the number of nodes in the hidden layer. Regarding the results of these tests, the neural network with 30 nodes in the hidden layer was selected as the best architecture for the model. The best results for each alternative are shown in Table 2.

Table 2. The training and testing RMS errors of three alternatives for the ANN architecture.

Number of Nodes in Hidden Layer	Training Error	Testing Error
25	0.323	0.371
30	0.107	0.181
35	0.329	0.362

After initial training of the neural network, testing the adaptation algorithm was done. Immediately upon running, the Data Comparison procedure started, it was seen that none of the current data pairs were identical to the data pairs in the initial data set. The current data pair was reflecting changed relationships between the inputs and outputs. Therefore one of the data pairs in the initial data set was replaced with the current data pair on each sampling cycle. The portion of new data increased rapidly but the error of the model was generally in the range of low to medium. The first retraining occurred after only 89 data pairs had been added to the initial data set. The training and testing errors after the first retraining process were respectively 0.104 and 0.443. The large testing error is likely due to a large diversity in the data pairs because of significant changes in the process.

This result is rather unsatisfactory since the "old" data pairs used for initial training were being removed from the data set indicating that the collected data and hence, the trained network were unreliable. The historical data were collected over one year prior to implementation (July-September 1997) and since then, there were many process changes that occurred in the plant such as:

- AG mill installation replacing rod mills.
- Difficulty in the reconciliation of tonnage rate to Line 3.
- Change in soda ash consumption measurement.
- Measurement errors in certain data input

The most important change was installation of the AG mill. AG milling has considerable effects on the variables and their relations such as grind size, ore liberation and pulp chemistry. The AG mill replaced all three rod mills (Line 1&2 and Line3), so the tonnage rate fed to the AG mill is much higher than that of the L3 rod mill. For purposes of applying the model, this tonnage was divided by three. However, it is possible that the tonnage to Line 3 circuit is different than this value from time to time. Thirdly, the soda ash consumption in L3 rod mill that was one of the important inputs into the initial neural network had a range of 1000-4000 g/t. This variable is now replaced by soda ash consumption in the AG mill with a range of 1000-7000 g/t. As well, some inputs such as the density and mass flowrate of the 1st Cu/Pb cleaner feed, the xanthate consumption (g/t) in the Cu/Pb rougher stage and the cell level measurements had many errors for the entire historical data collection period which were only recognized in hindsight.

The Novelty Filter in G2

During implementation of this project, we were asked on several occasions: "why develop such a complex, customized, adaptive algorithm when G2 already has such a procedure?" Well, there are very important differences in the logic used in the G2 adaptive algorithm and the one developed in this project.

In G2, if the Novelty Filter is attached to the data set collecting data pairs for training and testing, the data set is protected from being filled up with redundant data. Whenever a data pair is added to the data set, the Novelty Filter checks if there are more than a specified number of data pairs within a specified vicinity of the new data pair input variables and places such data pairs into a cell with the new data pair. The attributes of Points per Cell and Cell Size are configured by the user. An incoming data pair is judged to be novel if either of the following criteria is satisfied:

- If the cell contains only the newly received data pair, the data pair is novel.
- If the cell contains other data pairs, the filter averages the output values for those data pairs. If these averages differ from the new data output vector by a defined amount, the new data pair is novel and the other data pairs are discarded.

For adaptation, G2 suggests attaching a Novelty Filter to the data set. The Novelty Filter determines when the number of novel data pairs reaches a certain level, and sends a signal to the Training block to begin retraining. Now that both algorithms have been described the differences are clear:

- According to our algorithm, the frequency of a process relationship indicates its strength and is very important. Therefore in initial training, we did not try to collect only unique data pairs. Bear in mind that real-time data is read every 5 minutes allowing the creation of variety in the data as change occurs.
- In our adaptive algorithm, when data is replaced, although all data pairs that meet a particular condition are detected, only one is removed. We reason that if a certain relationship is identified several times, this indicates the importance of the relationship. Therefore it is incorrect to remove all data pairs in question. If a particular relationship no longer exists in the process, the arrival of new data pairs gradually removes the invalid ones and the data set is eventually adapted.
- The criterion to remove a data pair is also different. In G2, the age of the data pair is used but in our algorithm, the criterion is based on the Root Mean Squared Error which is preferred.
- Retraining in G2 depends on the amount of new data while our algorithm uses model error as well.

CONCLUSION AND RECOMMENDATIONS

A prototype adaptive Artificial Neural Network has been developed to predict the copper and lead assays of the Cu/Pb concentrate produced from the Line3 flotation circuit at the Brunswick Mining concentrator. The initial accuracy of the model was an RMS error of 0.107 for training data and 0.181 for testing data. Since the model actually predicts the rate of change in these assays and uses the actual assay from the previous cycle as input, these errors are viewed as acceptable starting points.

Following implementation, it was quickly realized that the initial training data set did not match well with current plant relationships because of significant circuit changes and significant error in the measurement of some input data. The model began adapting itself after 89 cycles (~7.5 hours). Following this adaptation, a similar training RMS error of 0.103 was achieved but the testing RMS error was now 0.443 which is probably due to extreme diversity among the data pairs in the data set. Despite such an apparently large error, trend plots show reasonable correlation with the actual assays. As additional retraining occurs, we anticipate the error will settle into the range of 0.10-0.20 and the retraining period will increase to ~7 days.

The need for phase lag calculations to synchronize data was apparent from the preliminary analysis of the data set. With real-time adaptation of these phase lags, we are able to adjust the model for changes in lag times on each specific variable. Future consideration should be given to separation of this phase lag into two components -- a pure dead time lag and a first order process time constant.

The model should be applied to other process assays. A separate neural network should be set up to predict each of the Cu, Pb and Zn assays of the Cu/Pb tailing, the Zn concentrate and the final tailing. As well, the Zn assay of the Cu/Pb concentrate may also be a useful variable to consider predicting. One key issue will be the cycle time required to run 12 networks in parallel with at least one network being retrained at the same time. We expect the time for these activities will be within the 5 minute cycle time of consecutive on-line XRF assays using current hardware and software. Recent data should be used to build these models as the data set is likely to be more representative of current ore and plant conditions than data one to two years old. Of the 58 input variables, about 10-12 were actually significant. The algorithm should be adapted to examine the variables and eliminate from training and running of the model, those which are unimportant.

In future work, by applying the established model and using the knowledge of expert operators and metallurgists, a compensating feed-forward model-based control can be developed and applied.

ACKNOWLEDGEMENT

The authors wish to express their appreciation to the workers at Brunswick Mining and Smelting Ltd., who assisted in this work through supply of the software and data.

REFERENCES

1. P.J. Poirier, H. Raabe, J.A. Meech, 1993. Using Froth Identification in an Advisory Expert System for Copper Flotation Operations. Proc. 25th Canadian Mineral Processors Conf., Ottawa, #36. pp.14.
2. Fletcher, R., 1980. Practical Methods of Optimization, Wiley and Sons, NY., 1, pp378.

Multivariable Predictive Neuronal Control Applied to Grinding Plants

Manuel Duarte M.*, Alejandro Suárez S. and Danilo Bassi*****

* D. Ing. Eléctrica, U. de Chile, Casilla 412-3, Santiago - Chile

** D. de Electrónica, U.T.F. Sta. María, Casilla 110V, Valparaíso - Chile

*** D. de Informática, U. de Santiago, Casilla, Santiago - Chile

ABSTRACT

This work investigates the use of a direct neural network predictive controller applied to a grinding plant. A phenomenological model of the grinding plant is used to simulate the control strategies. The model is based on a mass balance and power consumption of the mill containing 32 particle size intervals. The controller neural network is trained by using an estimation of the error. Several tests are performed driving the nonlinear process to an operation point and then controlling it by training the N.N. on line, which enables monitoring of the range over which the neural controller is still valid, without having to conceive a linear model of the process.

INTRODUCTION

Artificial neural networks (ANN) have been developed to approximate relationships between multiple inputs and outputs [1,2,21]. The neural network (N.N.) is trained to develop these relationships with pertinent data. ANNs have received considerable attention in multivariable control system applications [16,17,18].

The model predictive control has been used successfully in several industrial plants. In the majority of these applications a linear model is used to predict the behavior of the plant over a specific horizon of interest [12,13]. As the majority of real industrial processes exhibit a nonlinear characteristics, some researchers have extended the model predictive control technique in order to incorporate nonlinear models [14,15].

One of the most promising ways of avoiding the problem of defining the predictive model in an MPC is by using a N.N. as a black box representing the behavior of the nonlinear process [19,20]. Grinding control has been readily and thoroughly studied as grinding processes are costly and easily perturbed by changes in feed conditions.

Methods ranging from simple to complex control techniques have been used in the control of grinding circuits. One of the simplest forms is the use of SISO PID control loops, which is particularly useful if the loops are decoupled and PIDs are well tuned. However, since decoupling rarely occurs, multivariable techniques such as the Inverse Nyquist Array (INA) are much more appropriate [3, 4].

Modern control techniques including optimization are described by Herbst and Rajamani [5,6]. The method is based on a simple but accurate process model with an estimator such as the Kalman Filter used to predict model parameters and states. An optimization algorithm determines the necessary control actions.

An expert system employed in grinding and flotation control is described in [7]. The expert system based on pattern recognition is used as supervisory control which in turn is used to define the set points.

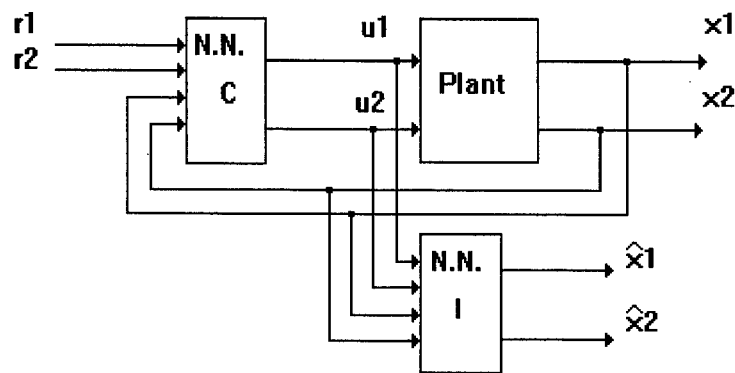
Autogenous or semi-autogenous systems are particularly difficult to control. The problems associated with these types of system are described by Duckworth and Lynch [8]. The difficulty is caused by the greater influence of feed non-uniformity supplied to the mill grinding circuit; if the feed stream does not contain sufficient average size particles then circuit behavior rapidly deteriorates. Changes in mineral hardness and size distribution in the feed rate can have even greater effects upon circuit behavior

In this paper the control strategy proposed in [11] is evaluated upon a multivariable dynamic model of the CODELCO-ANDINA grinding plant. This plant has three sections (A, B and C), each one of them formed by a bar mill and three ball mills.

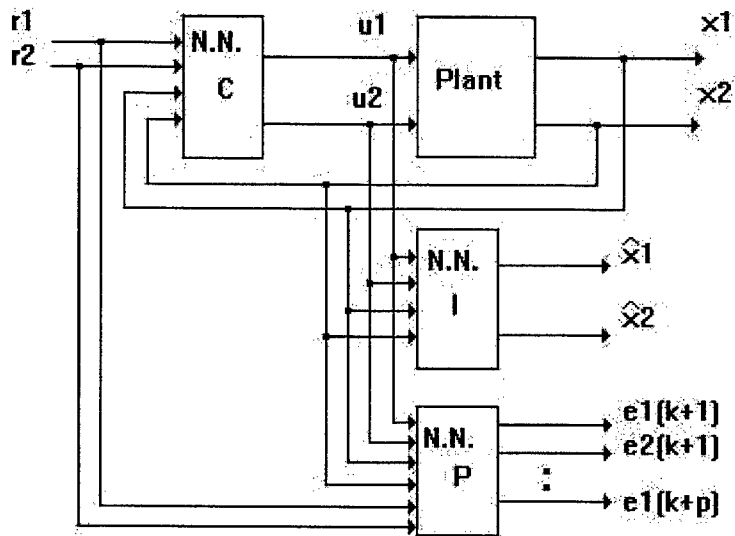
The proposed schemes, with and without prediction, are compared with the multivariable INA technique presented in [10] on the same plant. The phenomenological model used in the study is described in [9] while the predictive N.N. controller was developed in [11].

MULTIVARIABLE CONTROL SCHEME

Within these multivariable control configurations a comparison is made between a neural model with and a neural model without prediction error. The development of these schemes is presented in [11]. Figure 1 shows both MIMO schemes with the grinding plant's process variables.



a.



b.

Fig. 1. Neural control configuration with process model.
a. Without predictive network. b. With predictive network.

NEURAL NETWORK CHARACTERISTICS

It was possible to determine the delays to be considered in the ANN following a series of preliminary tests and studies. For the Dynamic Simulator it was determined that the best choice was that of nonlinear network structures with the following characteristics:

Controller Network

- Four layers (Input, two hidden and output)
- Input layer formed as:
 - Two references with four delays each.
 - $[r1(k-1), r1(k-2), r1(k-3), r1(k-4), r2(k-1), r2(k-2), r2(k-3), r2(k-4)]$
 - Two outputs feedback with four delays each.
 - $[x1(k-1), x1(k-2), x1(k-3), x1(k-4), x2(k-1), x2(k-2), x2(k-3), x2(k-4)]$
- First hidden layer with 10 neurons and activation function **tanh**.
- Second hidden layer with 5 neurons and activation function **tanh**.
- Output layer with 2 neurons and **linear** activation function.
- Sampling period 0.1 [min].

Identifier Network

- Four layers (Input, two hidden and output)
- Input layer formed as:
 - Two inputs with four delays each.
 - $[u1(k-1), u1(k-2), u1(k-3), u1(k-4), u2(k-1), u2(k-2), u2(k-3), u2(k-4)]$
 - Two outputs feedback with four delays each.
 - $[x1(k-1), x1(k-2), x1(k-3), x1(k-4), x2(k-1), x2(k-2), x2(k-3), x2(k-4)]$
- First hidden layer with 10 neurons and activation function **tanh**.
- Second hidden layer with 5 neurons and activation function **tanh**.
- Output layer with 2 neurons and **linear** activation function.
- Sampling period 0.1 [min].

Predictive Network

- Four layers (Input, two hidden and output)
- Input layer formed as:
 - Two inputs with four delays each
 - $[u1(k-1), u1(k-2), u1(k-3), u1(k-4), u2(k-1), u2(k-2), u2(k-3), u2(k-4)]$
 - Two outputs feedback with four delays each.
 - $[x1(k-1), x1(k-2), x1(k-3), x1(k-4), x2(k-1), x2(k-2), x2(k-3), x2(k-4)]$
 - Two references with four delays each.
 - $[r1(k-1), r1(k-2), r1(k-3), r1(k-4), r2(k-1), r2(k-2), r2(k-3), r2(k-4)]$
- First hidden layer with 10 neurons and activation function **tanh**.
- Second hidden layer with 5 neurons and activation function **tanh**.
- Output layer with 2 neurons and **linear** activation function.
- Sampling period 0.1 [min].
- Three prediction periods.
 - $[e1(k+1), e1(k+2), e1(k+3), e2(k+1), e2(k+2), e2(k+3)]$

GRINDING PLANT

The CODELCO-Andina grinding plant has three sections denoted A, B and C, each one consisting of a bar mill and three ball mills operating in a closed inverse circuit with their respective hydrocyclones. Section C (shown in Figure 2) was chosen for analysis simply because it has the best instrumentation.

The model used in this study corresponds to that implemented in a dynamic simulator presented in [9] the simulator was programmed in Turbo Pascal 7.0.

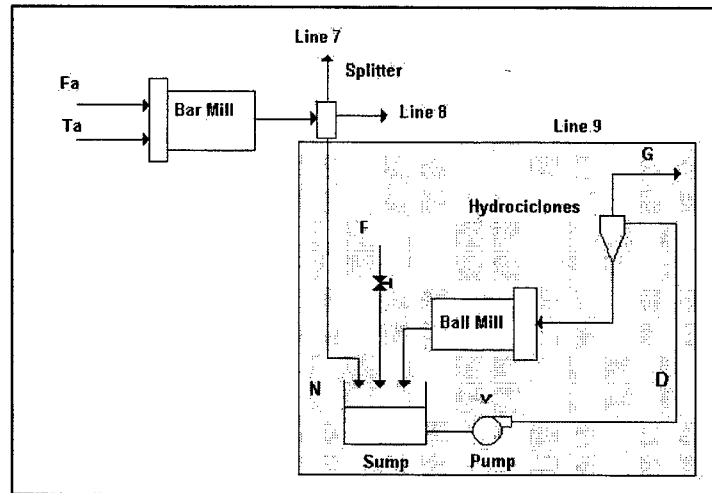
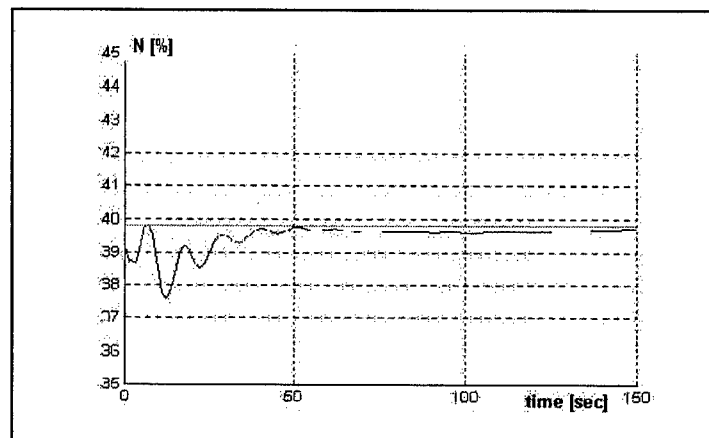


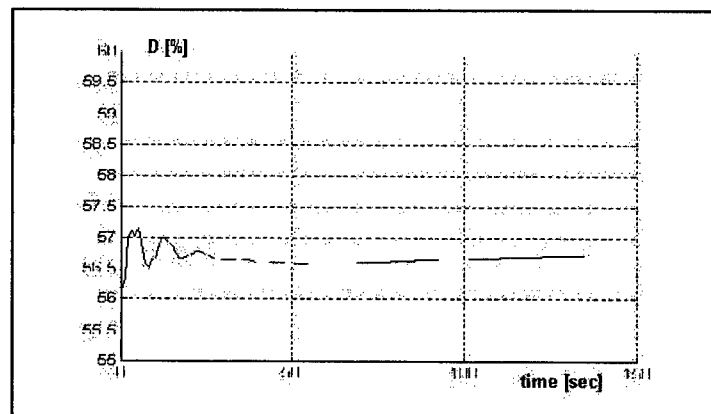
Fig. 2. Line 9 of Section C CODELCO-Andina grinding plant.

SIMULATION TESTS

The dynamic simulator was controlled using both predictive ANN schemes as described above. Fig. 3 and 4 illustrate the system behavior when control passes from manual to automatic. Transient adaptation is observed for a short period that is necessary to enable the ANN to stabilize around the new operating point.



a.



b.

Fig. 3. Output variables when control is changed from manual to automatic.

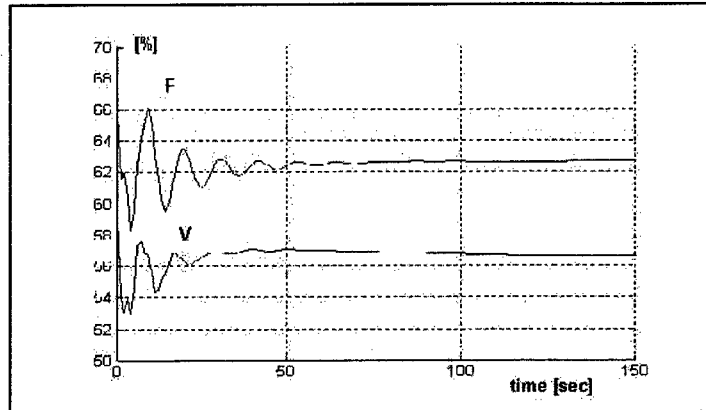
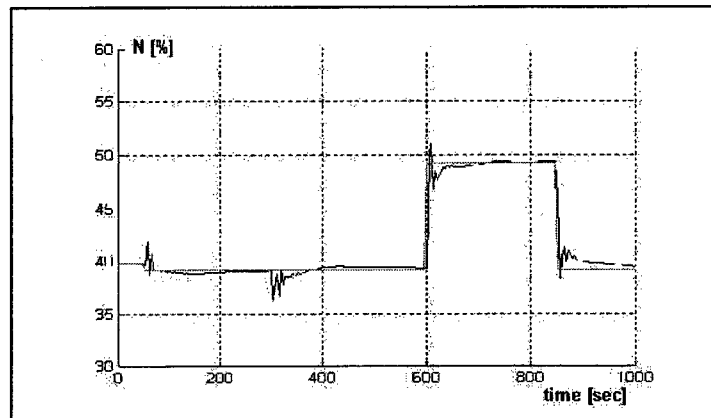
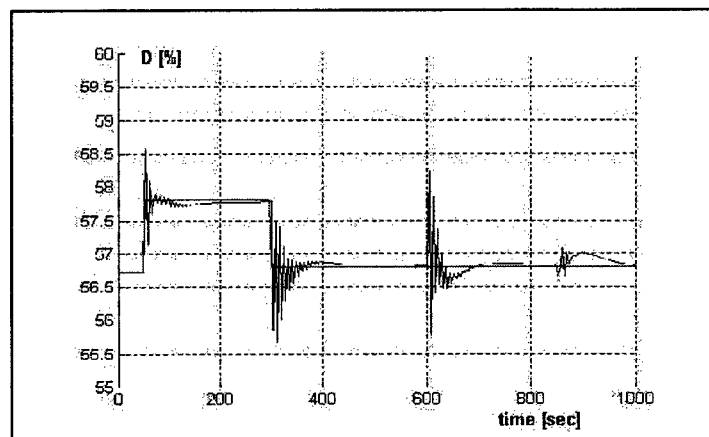


Fig. 4. Control variables when control is changed from manual to automatic.

Figures 5 and 6 show the response of the neural system without prediction. The evolution of variables N and D are plotted when changes in their references are produced. A small coupling between de D and N is observed in Figure 5, both during the transient and steady state regime. An oscillatory transient remains, which is more vigorous for density D.



a.



b.

Fig. 5. Neural control without prediction. a) Output variable N. b) Output variable D.

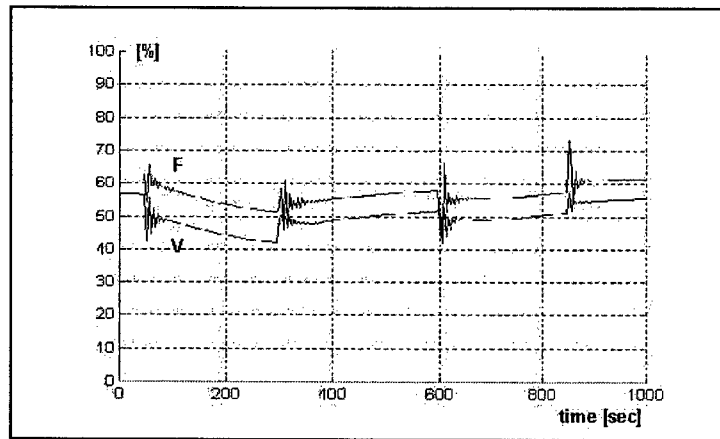
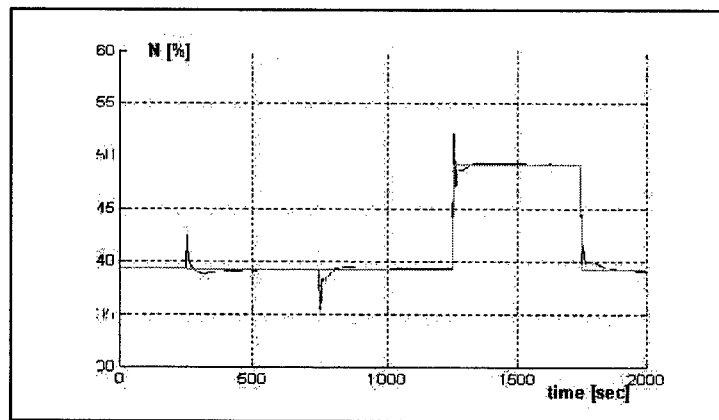
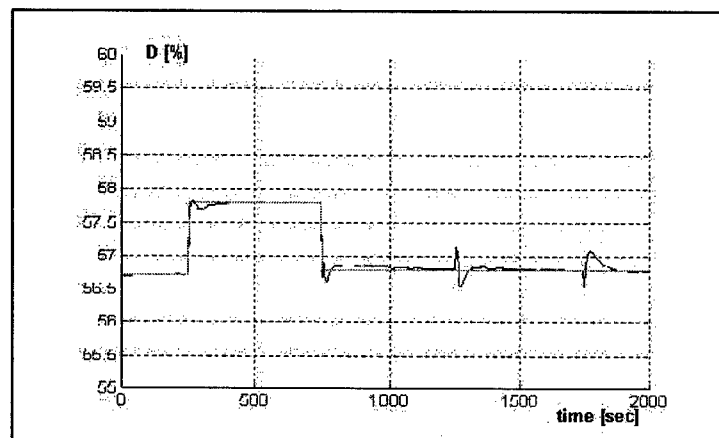


Fig. 6. Neural control without prediction. Control variables V and F.

Neural control with prediction is plotted in Figures 7 and 8. Plots for variables N and D are shown following changes in their references. Figure 8 reveals a smaller coupling effect than those shown in Figure 6, under both transient and steady state regimes. The oscillatory effects are greatly diminished with this control scheme when compared to Figure 6.



a.



b.

Fig. 7. Neural control with prediction effects. a. Output variable N. b. Output variable D.

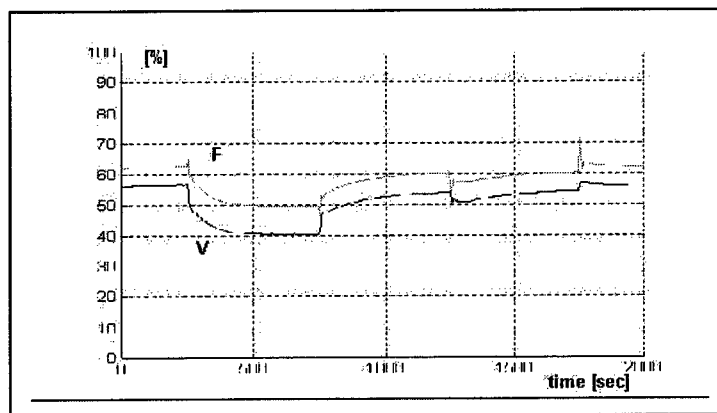


Fig. 8. Neural control with prediction effects. Control variables V and F.

CONCLUSION

Generally speaking, longer training periods are needed if non linear neural networks are used in control schemes due to the greater number of weights to be adjusted.

Insufficient training generates oscillations in the transient response. The output improved, in all tested cases, when predictive effects were included that diminished oscillations in the transient response.

Due to the length of time needed to train the networks, it is more effective to update the network's training on line at a low training rate. The only necessary off line training is then the initialization period.

The results presented in this work are clearly better than those reported in [10] since good decoupling is achieved between the output variables and, independently, a more rapid control action is obtained (results not shown).

ACKNOWLEDGEMENTS

The research presented in this paper has been funded by CONICY-Chile, through grant FONDECYT 1970351.

REFERENCES

1. Narendra, K.S., Parthasarathy, K., 1990. "Identification and control of dynamical systems using neural networks". IEEE Trans. Neural Net., 1, 4-27.
2. Narendra, K.S., Parthasarathy, K., 1991. "Gradient methods for optimization of dynamical systems containing neural networks". IEEE Trans. Neural Net., 2, 252-262.
3. Hulbert, D.G., Craig, I.K., Coetzee, M.L., Tudor D., 1990. "Multivariable Control of a Run-of-Mine Milling Circuit." J. South African. Inst. Min. Metall., 90(7), 173-181.
4. Jamsa-Jounela, S.L., 1990. "Modern Approaches to Control of Mineral Processing." Acta Polytechnica Scandinavica, Mathematics and Computer Science Series No. 57, Helsinki.
5. Herbst, J.A., Rajamani, K., 1982. "The Application of Modern Control Theory to Mineral-Processing Operations." Proc. 12th CMMI Congress, South African Inst. Min. Metall..
6. Rajamani, K., Herbst, J., 1991. "Optimal Control of a Ball Mill Grinding Circuit-II. Feedback and Optimal." Control. Chem. Eng. Science, 46 (3), 871-879.
7. Jamsa-Jounela, S.L., 1982. "Experiences in the Use of Expert System in Grinding and Flotation Process." Control at the Siilinjärvi Concentrator. XVII Inter. Min. Proc. Cong., CIM, Toronto, Canada.
8. Duckworth, G.A. and Lynch, A.J., 1982. "The Effect of Some Operating Variables on Autogenous Circuits and their Implication for Control." XVII Inter. Min. Proc. Cong., CIM, Toronto, Canada.

9. Duarte, M., et al., 1994. "Simulador Dinámico de una Sección de la Planta de Molienda de Codelco-Andina." Informe MC-16, U. de Chile, F. C. F. M., Dep't. de Ingeniería Eléctrica. Santiago, Chile.
10. Contreras, M., 1995. "Control Multivariable Frecuencial de una Planta de Molienda de Minerales. Trabajo de Tesis." U. de Chile, F. C. F. M., Dep't. de Ingeniería Eléctrica. Santiago, Chile.
11. Suárez, S., 1998. "Nueva Arquitectura de Control Predictivo para Sistemas Dinámicos No lineales usando Redes Neuronales". Tesis de Doctorado en Ciencia de la Ingeniería. U. de Chile. Dep't. Ing. Eléctrica, Santiago, Chile.
12. Cutler, C., Ramaker, B., 1980. "Dynamic Matrix Control: A computer control algorithm". Proc. 1980 Joint Automatic Control Conference.
13. Garcia, C.E., Morari, M., 1982. "Internal Model Control. A unifying review and some new results". Ind. Eng. Chem. Process Des. Dev. 21, 308-323.
14. Peterson, T., Hernandez, E., Arkun, Y., Schork, F.J., 1989. "Nonlinear predictive control of a semi batch polymerisation reactor by an extended DMC". Proc. 1989 Amer. Control Conference, 1534-1539.
15. D.D. Brengel and W.D. Seider, 1989. "Multistep nonlinear predictive control". Ind. Chem. Eng. Res. 28, 1812-1822.
16. Draeger, A., Engell, S., Ranke, H., 1995. "Model predictive control using neural networks". IEEE Control Systems Magazine, 61-66, Oct.
17. Murray-Smith, R., Sbarbaro, D., Neumerkel, D., 1992. "Neural networks for modeling and control of a nonlinear dynamic system." IEEE Symp. on Intelligent Control, Glasgow, Scotland. 122-127.
18. Haykin, S., 1994. "Neural networks a comprehensive foundation". Maxwell Macmillan Inter., NY.
19. Thibault, J., Grandjean, B.P.A., 1992. "Neural Networks in Process Control: a survey." Advanced Control of Chemical Processes. (K. Najim, E. Dufour, eds.), IFAC Symposium Series #8, 251-260.
20. MacMurray, J., Himmelblau, D., 1992. "Identification of a Packed Distillation Column for Control via Artificial Neural Networks". Proc. American Control Conference, San Francisco, CA, 1455-1459.
21. Hornik, K., Stincombe, M., White, H., 1990. "Universal Approximation of an Unknown Mapping and Its Derivatives Using Multilayer Feedforward Networks." Neural Networks 3, 211-223.

Practical Neural Network Applications in the Mining Industry

L. Miller-Tait, R. Pakalnis

Department of Mining and Mineral Process Engineering,
University of British Columbia,
Vancouver, B.C., Canada

ABSTRACT

The mining industry relies heavily upon empirical analysis for design and prediction. Neural networks are computer programs that use parallel processing, similar to the human brain, to analyze data for trends and correlation. Two practical neural network applications in the mining industry would be rockburst prediction and stope dilution estimates. This paper summarizes neural network data analysis results for a 1995 Goldcorp/Canmet study on rockbursting and a 1986 UBC/Canmet study on open stope dilution at the Ruttan Mine.

INTRODUCTION

Many aspects of mine design are based upon empirical data. Neural Networks analyze data and make predictions based on previous results. Neural networks have advantages over conventional empirical design approaches. These advantages include;

- Neural networks can easily use multiple inputs to analyze data,
- By using multiple hidden layers and nodes neural networks investigate the combined influence of inputs,
- Neural networks can be easily retrained as new data becomes available making them a more dynamic and flexible empirical estimation approach,
- Neural network software is inexpensive and easy to use,
- Neural networks have demonstrated a more accurate empirical estimate over conventional methods.

The advantages of using neural networks is illustrated in a rockburst prediction example and an open stope dilution example.

ROCKBURST PREDICTION

The first example of a potential situation where neural networks could be useful in the mining industry is the prediction of rockbursts through physical inputs. To quote directly from the Ontario Ministry of Labour "...we do not have the ability to predict when and where rockbursts will occur, and the experts in the field agree that we are not close to make such predictions" [1]. Between 1984 and 1993 eight underground miners were killed in Ontario due to rockbursts. This accounted for approximately 10% of underground fatalities during this period. If neural networks were to have success in predicting where rockbursts would occur additional ground support, remote equipment, and/or design modifications could reduce or possibly eliminate fatalities due to rockbursting. As safety is the primary responsibility of mining engineers, the potential for neural networks to assist in predicting rockburst inputs should be investigated.

In 1995, a joint project was completed by Goldcorp Inc. and Canmet called "Development of Empirical Design Techniques in Burst Prone Ground at A. W. White Mine" [2]. Part of the study was to collect input information on rockburst, caving, ground wedge, and roof fall failures at the A. W. White Mine between 1992 and 1995. This resulted in a failure database consisting of 88 ground failures with corresponding inputs for each failure. The six inputs collected for each failure were RMR [3], Q [4], span [5], SRF [2], RMR adjustment, and depth. These input factors were set up and run in a neural network with 73 examples being used for training and 15 examples being used to test the network. The output factor, stability, can be one of four failures [2] - PUN-RF (potentially unstable roof fall), PUN-GW (potentially unstable ground wedge), BUR (rockburst), and CAV (cave). A brief description of the input and output factors are listed below.

Input factors:

RMR - The RMR system, initially developed by Bieniawski in 1973, [3] bases rock mass quality on five parameters. These parameters are:

- uniaxial compressive strength of the rock
- rock quality designation (RQD)
- spacing of discontinuities
- condition of discontinuity
- ground water conditions.

These factors are given a numerical value and totalled together to get an RMR value. This value will be a number between 0 and 100 with zero being very poor rock and 100 being extremely good rock. The ground water conditions were assumed to be dry conditions.

Q - The Q factor refers to the rock quality tunnelling index [4]. Developed in 1974, by Barton, Lien and Lunde, from the Norwegian Geotechnical Institute, the Q factor is based on six factors, which are:

- RQD - rock quality designation
- Jn - joint set number
- Jr - joint roughness number
- Ja - joint alteration number
- Jw - joint water reduction factor
- SRF - stress reduction factor.

The actual Q formula is $Q = RQD/J_n \times J_r/J_a \times J_w/SRF$.

The Jw/SRF factor was assumed to be 1.0 for this study because dry conditions are assumed stress is factored through modelling and strain measurements. The Q factor ranges on a logarithmic scale ranging from 0.001 to 1,000 where 0.001 is extremely poor rock and 1,000 is virtually perfect rock.

Span [5] - the meaning of span refers to the width of an underground opening in plan view. Span can be determined through the largest diameter of a circle within an underground excavation.

SRF' [2] - refers to the adjusting of RMR values relative to stress ratios and previous history of ground conditions. It does not refer directly to SRF used in the calculation of Q. Stress criteria is based upon the ratio of induced stress over unconfined compressive strength (UCS) of the rock.

Output Factors

Burst - refers to a stope in which a rockburst has occurred. A rockburst is an instantaneous rock failure in or about an excavated area characterized/accompanied by a shock or tremor in the surrounding rock.

PUN-RF - refers to potentially unstable ground with respect to a roof fall. A stope is considered potentially unstable if any of the following conditions occur [2]:

- the opening may exhibit strong discontinuities having orientations that form potential wedges in the back
- extra ground support may have been installed to prevent a potential fall of ground
- instrumentation installed in the stope has recorded continuing movement of the stope back
- there may be an increased frequency of ground working or scaling.

PUN-GW - refers to a stope considered potentially unstable due to the likelihood of a ground wedge failure. This is a subset of PUN-RF collected separately to identify areas where jointing may result in wedge failures.

Cave - refers to when uncontrolled ground failures result in caving.

NEURAL NETWORK ANALYSIS

The above inputs and outputs were run on a neural network to see if a neural network could predict output results from the input data and also to see which inputs had the greatest effect on output prediction. A two layer network consisting of 13 nodes was run for 10105 cycles reaching a 1.69 percent error. Seventy three observations were used to train the network. The remaining 15 observations were used to test the network's predicting ability.

The results of the neural network shows that the network correctly predicted all outputs from the training data. The reason that this is not surprising is that the network used these 73 observations for prediction training. However, the neural network also predicted burst conditions on the test data which was new data for the neural network. The network appears to have trouble distinguishing between PUN-GW and PUN-RF but predicted burst conditions on every occasion. The fact that burst conditions were predicted on each occasion was promising with respect to the possibility that neural networks may be a useful tool to predict rockbursts.

It appears from this database, that SRF has the most significant effect on predicting rockbursts. The bias node, Q, and adjusted RMR are also significant while RMR, span, and depth appear to have a lesser effect. It is not surprising that SRF has the most significance as it is a factor given to rock according to its previous history of burst proneness. A larger database with stable openings included is necessary to gain confidence in neural network predictions and the influence of input factors.

This example, by simultaneously analyzing several inputs, shows that neural networks can provide an effective tool for predicting rockbursts. Further work varying error, the number of nodes, layers and cycles could improve the network using this database. However, a larger database with more input factors could make a neural network a more effective burst predicting tool that could be practically applied in the mining industry.

NEURAL NET/FORMULA DILUTION PREDICTION COMPARISON

In an effort to compare neural net results with conventional formulae estimates, neural net predictions were compared with three formulas developed through a database collected from the Ruttan Mine. The formulas being compared are from three slope configurations: isolated slopes, echelon slopes, and rib slopes [6].

Isolated Slopes (61 observations)

1. $\text{Dil}(\%) = 5.9 - 0.08(\text{RMR}) - 0.010(\text{ER}) + 0.98(\text{HR})$
Echelon Slopes (44 obs)
2. $\text{Dil}(\%) = 8.8 - 0.12(\text{RMR}) - 0.18(\text{ER}) + 0.8(\text{HR})$
Rib Slopes (28 obs)
3. $\text{Dil}(\%) = 16.1 - 0.22(\text{RMR}) - 0.11(\text{ER}) + 0.9(\text{HR})$
where:

DIL(%) - Slope Dilution (%), ie. 10%, DIL(%) = 10

RMR - CSIR Rock Mass Rating (%), ie. 60%, RMR = 60

ER - Exposure Rate as Volume removed (metres cubed)/mth/slope width (m)

HR - Hydraulic Radius (m) of exposed slope wall

Neural nets were developed from the same databases that these formulae were developed. The neural net predictions on unseen data (not in the original databases) were compared with the formulae estimates. This was done to provide insight into the effectiveness of neural net predictions compared with statistically developed formulae estimates. The neural networks were not optimized.

The neural networks were trained using one hidden layer of two unseen nodes. Each neural net was trained on the entire original database for each stope configuration. The neural net trained on the rib stope database was trained to eight percent error, the neural nets trained on the echelon and isolated stope databases were each trained to 10 percent training error. These neural nets were then used to predict dilution on new unseen data and compared with the respective formula dilution estimates. The differences in the actual dilution from the neural net and formula predictions were compared for each stope and the combined averages of the stopes for each stope configuration. Figure 1 charts the average percent error between the neural net and formula predictions.

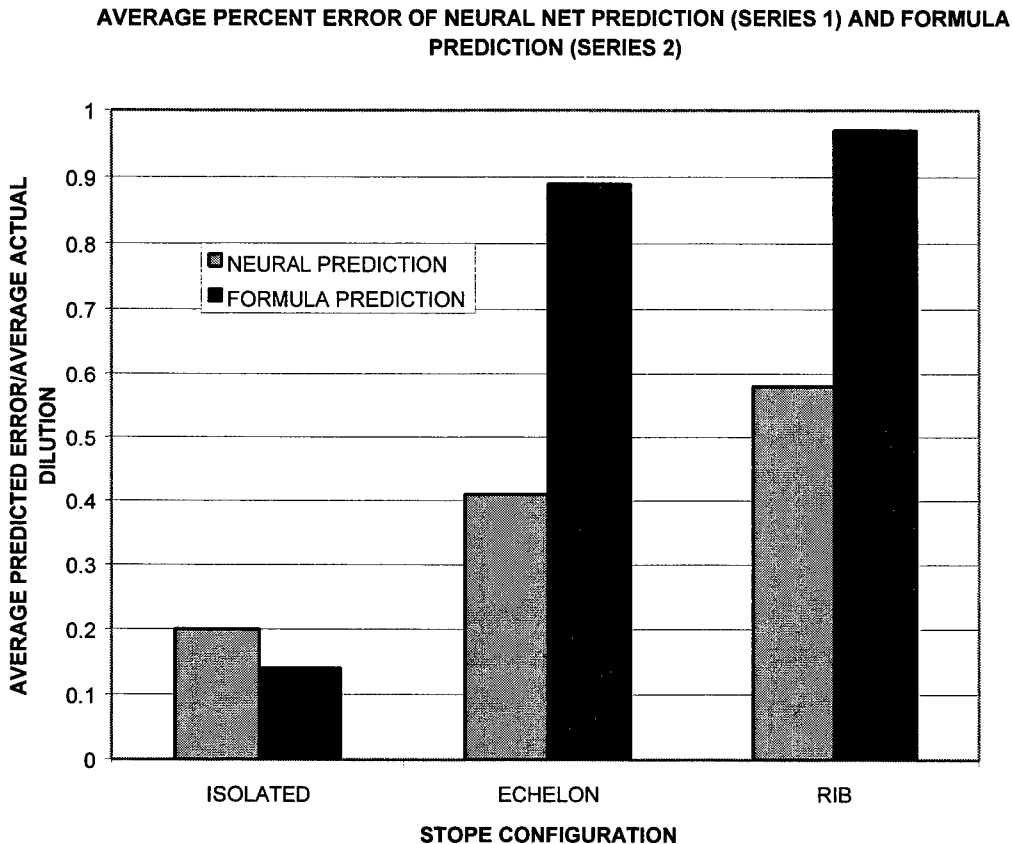


Fig. 1. Average Neural Net / Formula Error Over Actual Average Dilution

For the unseen rib stope data, the neural net had an average error of 3.2 percent dilution and the formula had an average error of 5.1 percent dilution. For the echelon unseen stope data the neural net had an average error of 1.8 percent while the formula had an average error of 3.9 percent. For the two unseen isolated stopes the neural net had an average error of 0.9 percent while the formula had an average error of 0.6 percent. As the rib stope and echelon unseen databases were significantly larger than the isolated stope unseen database the neural network showed a clear improvement over the formula estimates. The improved performance of the neural net predictions in this example over the statistically derived formulas suggest that neural nets can have better predictions than conventional formula predictions.

CONCLUSION

The Goldcorp/Canmet example shows that neural networks can provide an effective tool for predicting rockbursts. Further work varying error, the number of nodes, layers and cycles could improve the network

using this database. However, a larger database with more input factors would make a neural network a more effective burst predicting tool. Additional inputs for each failure may include: induced stress (map3d), hydraulic radius, presence of raises, microseismic data, faults or dikes, ground support; type of heading, and if active mining is in the vicinity.

The improved accuracy of the neural net dilution predictions over the statistically derived dilution formulas suggest that neural nets can produce more accurate estimates over conventional empirical methods. The need to having adequate amounts of input data was demonstrated as the training error improved on the smaller stope dilution databases but the accuracy of predictions on unseen data decreased for the smaller databases. Besides improved accuracy, the added neural net advantages of multiple inputs and the continuous ability to retrain the neural nets should improve empirical estimates in the mining industry.

REFERENCES

1. Caron, M., 1995, Ontario Ministry of Labour Rockburst Recommendations, Ontario Ministry of Labour.
2. Mah, P., 1995, Development of Empirical Design Techniques in Burst Prone Ground at A.W. White Mine; Canmet Project No.: 1-9180.
3. Bieniawski, Z. T., 1989, Engineering Rock Mass Classifications, New York; John Wiley & Sons.
4. Barton, Lien, Lunde, 1974, Classification of Rock Masses for the Design of Tunnel Support, Rock Mechanics Vol. 6, No. 4, 7 pp.
5. Lang, B., Pakalnis, R., Vongpaisal, S., 1991, Span Design in Wide Cut and Fill Stopes at Detour Lake Mine, 93rd AGM - CIMM, paper # 142, Vancouver.
6. Pakalnis, R., 1986, Empirical Stope Design at the Ruttan Mine, Sherritt Gordon Mines Ltd., University of British Columbia, Canada, 276 pp.

Neural Network-Based Resistance Spot Welding Control and Quality Prediction

Nenad Ivezic, John D. Allen Jr., and Thomas Zacharia

Computer Science and Mathematics Division
Oak Ridge National Laboratory
Oak Ridge, TN 37831, USA

ABSTRACT

This paper describes the development and evaluation of neural network-based systems for industrial resistance spot welding process control and weld quality assessment. The developed systems utilize recurrent neural networks for process control and both recurrent networks and static networks for quality prediction. The first section describes a system capable of both welding process control and real-time weld quality assessment. The second describes the development and evaluation of a static neural network-based weld quality assessment system that relied on experimental design to limit the influence of environmental variability. Relevant data analysis methods are also discussed. The weld classifier resulting from the analysis successfully balances predictive power and simplicity of interpretation. The results presented for both systems demonstrate clearly that neural networks can be employed to address two significant problems common to the resistance spot welding industry, control of the process itself, and non-destructive determination of resulting weld quality.

INTRODUCTION

Several factors may influence the quality of a forming resistance spot weld. Among the more important of these are failures in weld tip geometry, improper alignment of welder electrodes and metal surfaces to be joined, dirt and corrosion on the electrodes and/or metal surfaces, and uncompensated variations in AC supply voltage. Each of these may influence the principal variables (the welder output current and weld tip voltage drop) to which the weld control/quality assessment system is permitted to have access. Of primary importance is the signature defined by the temporal variations of these two quantities.

The first task of the work reported here was to develop a system able to control the resistance spot welding process in real time. Realization of such a capability depended critically on developing methods to detect and compensate for influential factors on a short-enough time scale to support modulation of an evolving weld. The recurrent neural network system developed for this purpose is discussed in the next section.

The second task was to develop a system able to perform an a-posteriori weld quality assessment. In one approach, the recurrent neural network developed as part of the first task was adapted to produce evolving evaluations of welding signatures and extract from them a measure of weld quality. In an alternative approach, a static neural network was incorporated into an adaptive system that proved able to determine weld quality by predicting key characteristics of the weld – nugget size and indentation – and by the subsequent mapping of these characteristics onto a pass/fail classification.

EXPERIMENTS IN WELD EVOLUTION CAPTURE FOR PROCESS CONTROL AND QUALITY ASSESSMENT

Critically important to the production of quality spot welds is application of appropriate welding current for a time sufficient to ensure formation of the weld nugget and short enough to avoid undesirable effects (e.g., excessive denting of the material surfaces and thinning of the weld region due to prolonged application of electrode pressure). A reliable indicator of the onset of these and other problems is expulsion of weld material occurring when the forming weld nugget begins to exceed in volume that which can be effectively

contained within the electrode boundaries. Welding current should be (indeed, should already have been) removed when this onset is detected. In many spot welding control systems, this detection is effected by real-time analysis of derivatives of the weld resistance versus time history of the welding process. Since a control function that results from such an analysis, must occur after the analysis has determined the onset of expulsion, weld integrity may be irreversibly compromised before the welding process can be terminated.

Real-Time Weld Evolution Prediction

In preparing for the investigation of real-time prediction of weld characteristics, a commercial automotive spot welding apparatus and controller were installed in our laboratory and provided with a high-speed data acquisition system capable of acquiring weld data at speeds adequate for the purposes of developing a real-time process control system. Also included in the system was a "compensation coil" from which data could be derived for counteracting the inductive effects of the very large welding currents. Compensation data were also employed to support a sensing system that eliminated the necessity for making tip voltage measurements at the welding electrodes themselves, an important consideration in an industrial setting.

The first experiments with the laboratory welding apparatus were designed to demonstrate the feasibility of real-time prediction of dynamic resistance values during weld production. If such prediction was possible and reliable, these results could be trivially incorporated into a real-time weld-termination system and, with somewhat more effort, into a dynamic weld-control system as well. The network(s) employed for this purpose were simple feed-forward recurrent perceptrons enhanced with several sets of integrating shift-register input nodes through which signals fed back from an intermediate layer were cycled with each presentation at the normal input nodes of new welding data. Each set of shift-register nodes was characterized by a unique time constant, the individual values chosen to optimize the "historical" memory of the network. For these studies, the interval between measured welding data and the ensuing predicted values was limited to no more than three weld cycles. Figure 1 depicts the network predictions for weld resistance together with the observed resistance values for a series of 49 resistance spot welds of duration spanning a few to a few tens of weld cycles. Note that the welds were not obtained in succession, as the figure suggests. Rather, prediction results from the 49 welds were combined in a single data set to simplify subsequent analysis. In any case, the predicted values are so nearly identical to the subsequently observed values (which, for display and comparison purposes, were shifted into registry with the observed data) that the disparities are scarcely visible at the scale of the figure. Thus, it was concluded that, at the very least, dynamic control of the weld termination process could be effected with the developed system.

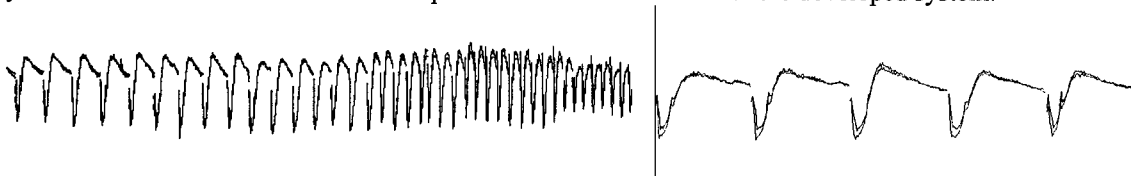


Fig. 1. a. Network prediction vs. observed resistance; b. Expanded representation of Fig. 1. a.

WELD QUALITY ASSESSMENT WITH ATTRACTOR NEURAL NETWORKS

Development of an attractor-based weld classification system represents a natural extension of the real-time dynamic weld-termination system described above. The method is suitable for application in the context of real-time welding environments and in a-posteriori assessment of previously acquired welding data.

Central to the attractor method is a novel dynamic network training mechanism derived from consideration of two aspects of weld classification. It would be entirely unreasonable to expect a system to determine weld quality on the basis of data representing only the first few weld cycles (or the middle or last few, for that matter). On the other hand, it is reasonable to expect the weld quality estimate to improve as more and more cycles of a particular weld event are completed. So whatever a system might do in response to temporally-evolving data for which a final "Good / Bad" rating is required, it should be done in a manner such that (1) very little response is given during the initial phase of data presentation (i.e., the first few weld cycles) and (2) a very stable response is obtained during the final moments of data presentation (i.e., the last few weld cycles). This requirement is realized in an attractor network by imposing on the training goal, a signal having the form of a (suitably offset) square root of a half-sine, a function with relatively modest and slowly varying

slopes at the two extrema and a much higher slope throughout a broad central region. During training, welds known to be good were forced to define a training function beginning at half full node output and following the square root of a half-sine function to maximum output at the end of the weld (regardless of the number of weld half-cycles in the complete weld). Welds known to be bad were forced to define a training function beginning at half full node output and following the square root of an inverted half-sine function to minimum output at the end of the weld (also irrespective of the total number of weld half-cycles). The endpoints for "good" and "bad" welds define two attractor basins for the network.

The network developed for the time-dependent classification task was very similar in form to the recurrent network described above and comprised an input layer, a single hidden layer, and an output layer. The input layer for a typical classification network included sixteen nodes devoted to the single-cycle near-peak values of the time-varying V and I (eight for V , eight for I), several nodes for the "static" weld parameters, and three node groups arranged as a sort of block shift register, each of the blocks being characterized by a successively longer integration time than its neighbor (and its signal source). "Previous-cycle" hidden layer outputs were mapped onto the nodes of the first of these three groups and subsequently shifted leftward one block per weld half-cycle. This network proved surprisingly adept at classifying weld quality as long as training data and test data were acquired with the same welding apparatus. Moreover, such a network on several occasions correctly assessed a weld that a human weld assessor had misclassified. Perhaps most intriguing is the fact that the attractor networks demonstrated a capability for recovering from initial error. Thus, the output response for a weld ultimately properly rated as "good" could, if welding conditions were initially ambiguous, swing in the opposite direction for several cycles before moving toward, and finally to, the proper classification endpoint (attractor). Those welds which properly fell along the boundary between good and bad led to oscillatory network responses with network output never moving far from half full output value.

Figure 2 shows trained network responses for a series of 30 welds of various durations. Of these, 10% were previously unseen by the network. Although overall classification accuracy was of the same order as that exhibited by the static network in the next section (95-98 %), the ability conferred by the attractor network to deal with dynamic systems would appear to provide advantage when "real-time" operation is required.

WELD QUALITY ASSESSMENT WITH STATIC NEURAL NETWORKS

The development of the static network method for weld quality assessment arose as a result of an attempt to develop software that would exhibit minimum sensitivity to the effects of significant, but not quantified, variability of welding process parameters. This section covers data analysis, neural network design, and results obtained in using the static neural network classifier. Two experimental designs, the second narrower in focus than the first, define the scope of the presented material. Figure 3 is a high-level representation of the components comprising the static weld classification system.

Preliminary Data Analysis and Neural Network Design

Data analysis for the static network system was performed in the context of three time-varying parameters, current through the secondary winding of the welder transformer (I), welder transformer secondary voltage (V), and the voltage appearing at the previously noted "compensation" coil (V_c). An exploratory data analysis performed on a number of other process parameters (e.g., SCR voltage, SCR firing angle, power factor) suggested that the information content of these variables was not significant.

In order to represent the time dependent I , V , and V_c values to the static network, these data were unfolded into the components of a static, multi-block, input vector of which each block represented data for a single weld half cycle. Desired outputs of the network classifier were taken to represent "PASS" and "FAIL" values for weld quality. It should be noted, however, that the network was not forced to learn the mapping of weld data to weld quality directly. Instead, the classifier system was constructed so its two outputs computed the values of two spot weld attributes on which classification is traditionally based, the nugget size (NS) and the indentation value (Ind). The computed values for NS and Ind were used, not directly in a classification formula, but to identify "updated" threshold values in the weld classification algorithm.

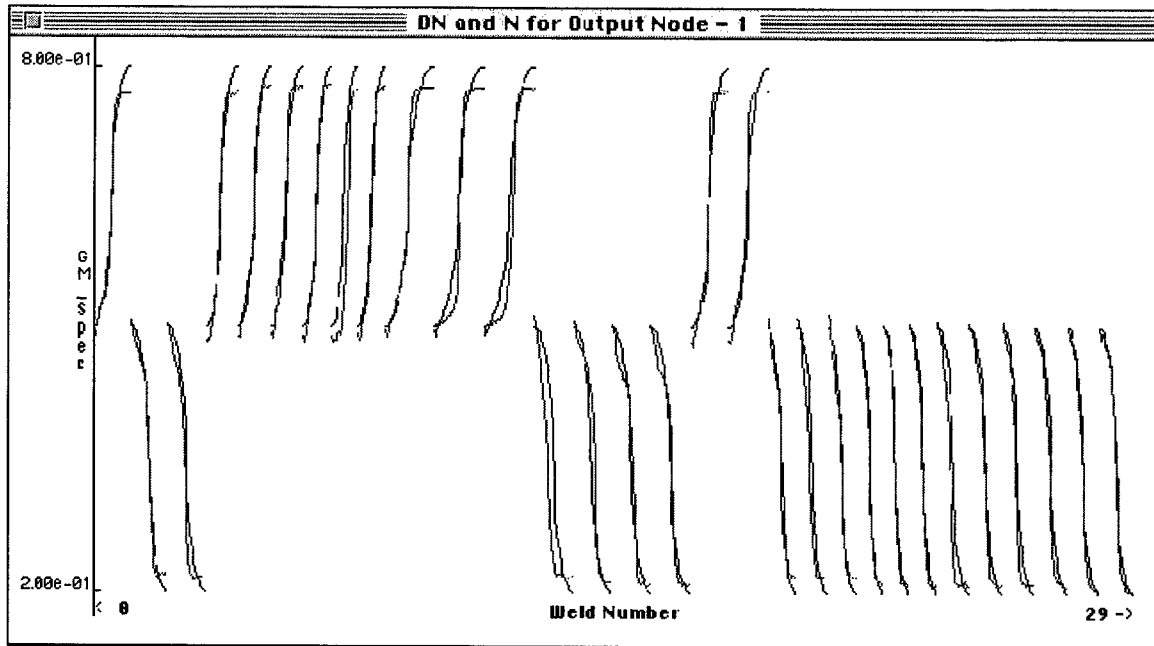


Fig. 2. Weld quality prediction with recurrent network.

The rationale for this treatment of network output followed from the requirement to manage the variations in unobserved parameters in an optimal manner with respect to the classification task. This requirement arose in turn from the earlier established fact that the information content of the measured variables (I , V , and V_c) did not appear sufficient to lead to construction of a static network-based classifier of the desired accuracy. This forced a trade-off among classifier characteristics. In particular, it was deemed better to be able to classify correctly all failing welds at the expense of incorrectly classifying some of the good welds as failing rather than to risk mis-identification of failing welds. Network-computed values for NS and Ind supported a classification threshold update scheme by which this trade-off was effected.

Experimental Design #1

The defining condition for the first experimental design was to maintain a realistic level of difficulty for the classification task in the context of a reduced, yet significant, variation in welding schedules. Controlled variables were limited to Tip Force, Number of Weld Cycles, Weld Current, and Upslope Current. A training data set of 144 experimental welds divided into three groups was created. The groups were defined by Tip Force values of 400, 500, and 600 lb., respectively. Within each of these three groups, the number of weld cycles was limited to the values 4, 6, and 8. Finally, for each fixed number of weld cycles within each of the 3 groups, both Upslope Current and Weld Current were varied simultaneously to obtain 16 different combinations. Sixteen additional welds were created for use as testing data. Parameters for this group of welds were intentionally chosen as duplicates of those of members of a previous subgroup.

Analysis of Experimental Results #1

The primary purpose of experiment #1 data analysis was to identify a "signature" of the welding process that would allow reliable classification of resistance spot welds. Two important points emerged: (1) most process parameters (e.g., Power Factor, Firing Angle) do not appear to contain discriminating information relevant to the classification task; and (2) the raw, unprocessed input signals (V , I , and V_c) are not by themselves sufficiently discriminating for the static network-based classification task.

At this point, we made use of previous work on the Dynamic Resistance Curve [1] to try to find a more suitable representation for classification. The important result was that raw input signals (for V and I), when transformed into points on the dynamic resistance curve clearly contained discriminating information. The nature of signal acquisition used in these studies suggests that a compensation representation shown in Eq. 1 can recover the approximate form of the dynamic resistance for a $k = 5.0$.

$$R = (V - k \cdot V_c) / I$$

1.

Neural Network Classifier Design for Experimental Results #1

The architecture for the first static network-based classifier was a simple Perceptron comprising an input layer, a single hidden layer, and an output layer. The number of input units depended on the number of weld cycles in the experimental group with which the network was associated. For example, the second group of weld experiments was defined by 6 Up Slope cycles and 6 Weld cycles. So the associated network included 24 input units (two for each complete AC cycle, or one per each so-called "weld half-cycle"). The signal represented by each unit was taken as the value of the "dynamic resistance" computed according to Eq.1. Two output units were included in all static classification networks, one to compute the Nugget Size, the other for the Indentation value. Network training involved a form of cross-validation wherein the complete data set for any one experimental group was divided into various training and testing sets.

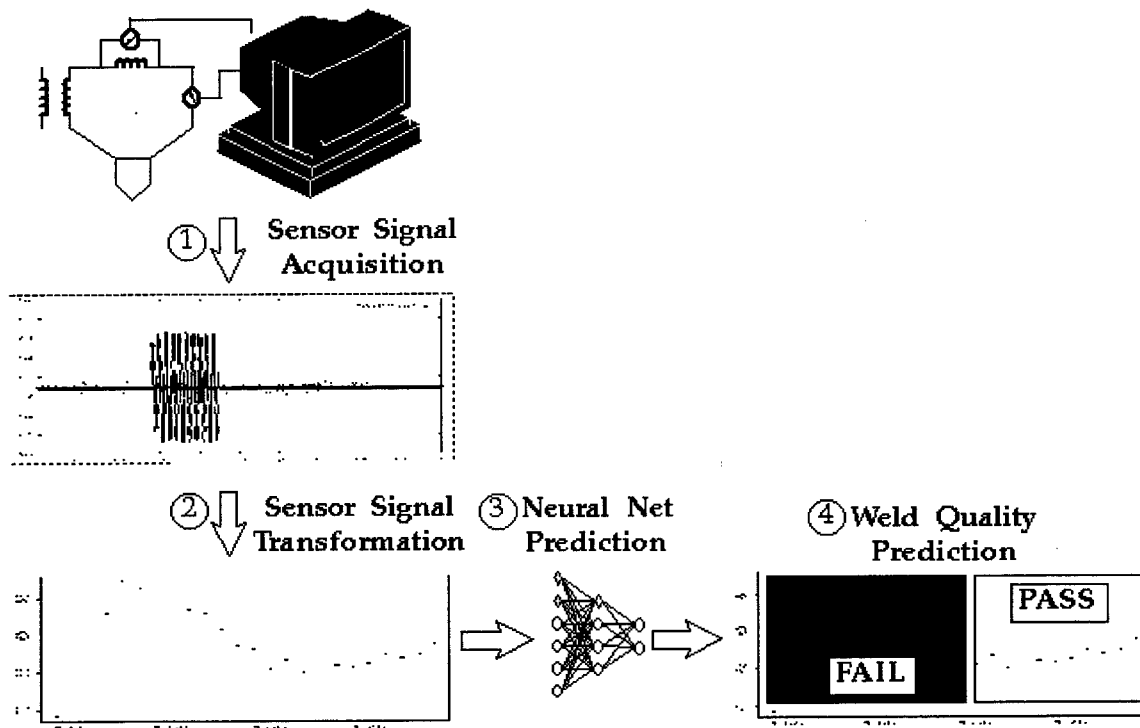


Fig. 3. Functional representation of "static" weld quality prediction system.

Although the classification accuracy observed in the cross-validation tests was observed to be in excess of 80%, the performance of the classifier was relatively poor in those cases for which the welds were either marginally passing or marginally failing. The experiment described in the next section was designed to investigate whether classification performance could be improved for such marginal cases.

Experimental Design #2

For the second experiment, the number of controlled variables whose values were varied was further constrained. The intent was to investigate whether classifier accuracy would increase given training cases covering a limited sub-space of the welding process parameters.

Data for the second experiment comprised 96 training instances divided into two groups and 16 testing instances. The Tip Force was held constant at 500 lb. and the Welding Cycle value was fixed at 6. The only variable parameters were the Upslope Current and the Welding Current. For the training set, these were varied to obtain 16 combinations. For each of these combinations, three experimental welds were generated so as to capture the effect of the unobservable variables. Each of the first two groups of experimental welds

thus bore $16 \times 3 = 48$ welds. The third testing group included a single weld for each of the 16 Upslope Current/Welding Current combinations.

Analysis of Experimental Results #2

Data analysis for the second experiment confirmed the previous finding that there is a significant correlation between the data points transformed so as to represent a dynamic resistance curve and the changes in Indentation and Nugget Size variables. However, it was also clear that, in borderline cases in which welds were either marginally passing or marginally failing, the constructed dynamic resistance curves did not exhibit visible discriminating characteristics. The task was now to determine if a neural network classifier could extract features that were not obvious from visual inspection of these data points.

Neural Network Classifier Design for Experimental Results #2

The neural network architecture for the second set of experiments was identical to that employed in the previous stage. However, a "power" representation of weld parameters (Eq. 2, below), rather than the expression of Eq. 1, was used to represent the dynamic resistance curve points to the network classifier.

$$P = (V - k \cdot V_c) \cdot I \quad 2.$$

Results of network training involving the power representation of dynamic resistance were substantially improved relative to those of the earlier employed "static" schemes. In the cross-validation tests, the classification performance was now observed to be in excess of 95%. When the constructed network classifier was tested on the cases from the previous experimental design #1, the same result was observed. This held true irrespective of the fact that the Tip Force in experimental run #1 was in some cases different from that employed in the present experiment. Parallel validation experiments were performed in the laboratory and welding shop environments with similar results.

CONCLUDING REMARKS

The studies reported here offer clear evidence that relatively simply obtained resistance spot welding data embody signatures that can be extracted and employed for weld quality assessment. The observed 95% classification accuracy for a series of experimental runs performed at significantly different times and for different process parameters suggests considerable potential for the static network method. The dynamic method, although less extensively studied, appears to offer considerable potential in applications for which the ability to perform classifications during the evolution of a physical process is of importance.

At the most general level, we assert that we have demonstrated that neural networks can be used to address two significant problems confronting the resistance spot welding industry: control of the process itself, and both real-time and rapid a-posteriori determination of completed weld quality. The primary objective of future work is to conduct experiments and modeling studies to validate neural network technology and determine its suitability for quality assessment within a broader industrial context. The main task to complete is to evaluate the sensitivity and robustness of the developed neural network-based approach.

ACKNOWLEDGEMENTS

This research was sponsored by the Office of Basic Energy Sciences, Division of Materials Science, U.S. Department of Energy, under contract DE-AC05-96OR22464 with Lockheed Martin Energy Research Corporation. This manuscript has been authored by a contractor of the U.S. Government under Contract DE-AC05-96OR22464. As such, the U.S. Government retains a non-exclusive, royalty-free license to publish or reproduce this contribution, or allow others to do so, for U.S. Government purposes.

REFERENCES

1. S.A. Gedeon, C.D. Sorensen, K.T. Ulrich, T.W. Eagar, 1987. Measurement of Dynamic Electrical and Mechanical Properties of Resistance Spot Welds. *Welding Journal*, December, 378-385.

Verifying Detected Facial Parts By Multidirectional Associative Memory

M. Kitabata, Y. Takefuji

Graduate School of Media and Governance, Keio University, Kanagawa, Japan

ABSTRACT

In this paper, we propose a neural network system for verifying whether a mouth or eyes can be extracted from an image area by Back Propagation (BP). It is necessary to test the proposed system in a noisy environment. In this paper, the model of neural network system for recognizing a mouth is based on the function of peripheral vision. In our research, a mouth has distinct properties of brightness in the right corner of the mouth, the left corner of the mouth, the tip of nose, and the nostril. Furthermore we discovered that humans commonly observe these properties of the mouth regardless of the brightness of lighting, different colors of the mouth, or different form of the mouth. By using these features, we designed an associative memory neural network for the verification.

INTRODUCTION

The face is one of the most important clues to identify a person in a scene in the visual system. In particular, an eye and a mouth are the important clues to recognize a person. In our research group, we designed a human recognition system using neural networks by detecting and grouping facial parts [1,2]. The system of detecting facial parts extracts an eye and mouth from an input image by BP. Although the system can extract eyes and mouth in short computational time, it does not perform well in noisy conditions. It is also difficult to normalize a position and a size of image. This detecting system often misidentifies the facial parts due to the positional gap or the different size of facial parts.

Since a human cannot completely recognize an object in a small image, a larger image to recognize it is often substituted [3,4]. The large size of a visual field to perceive an object is divided into some regions. Generally the visual information is processed on the center of the field with high resolution, and on the periphery of the field with low resolution [5]. The color is perceived by cones on the center of the retina. Humans do not distinguish color in the peripheral area of the visual field. Only the contrast of brightness is perceived by the rods on the periphery of the retina. In the 1930's, physiologists revealed that humans segregate a figure and ground by the contrast of brightness to perceive a figure in a scene. Human eyes can perceive an object in different brightness of lighting, by extracting the contrast between the figure and the ground. For example, we can read a book in a brightly or darkly lit room. This phenomena is called "light and dark adaptation of the eye" [6,7].

The purpose of our research is to simulate a biological model which can precisely recognize and verify an eye and a mouth extracted by BP. A system to verify an eye using an associative memory MAM has been proposed [8]. In this paper, we show the system performance. In the simulation, we tested image sizes between 16×16 and 52×52 , and found that 32×32 is the best size for our system to verify a mouth in an image. Our system recognizes a mouth by linking related features around the mouth even when the image contains noisy data.

EXTRACTING THE FEATURES OF THE MOUTH

In our research on recognizing the mouth in the image, we discovered that both corners of the mouth and the nose have some distinct properties as shown in Fig. 1. The ratio of the brightness was calculated on each area which showed the property of the mouth. These properties are observed in the closed or opened mouth, regardless of lighting conditions, individual variation, or the mouth form. This valuable information enables us to extract the mouth in an image.

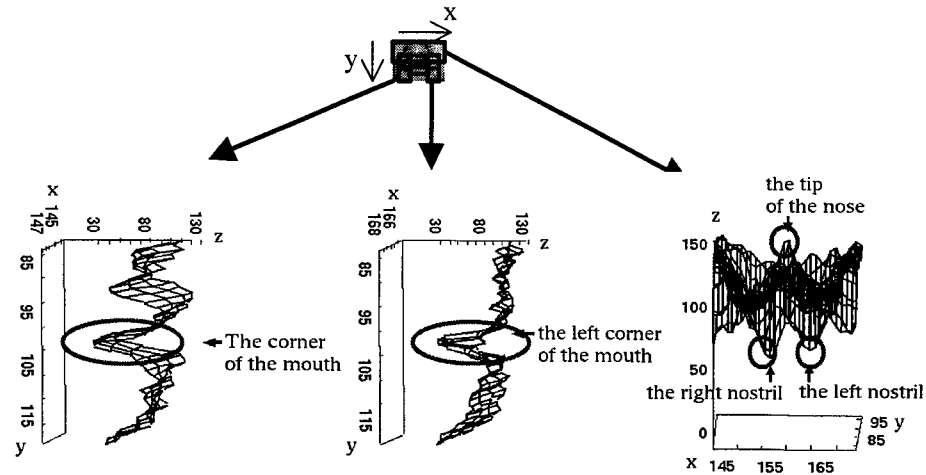


Fig. 1. The properties of brightness contrast are shown in the circled area with broken line: The xy-plane represents the pixel position in an image, and the z-axis represents the pixel brightness.

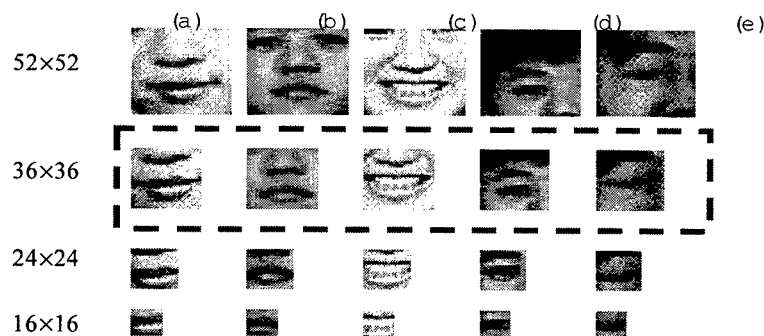


Fig. 2. The mouths and eyes in the small images: (a), (b) are the closed mouths, (c) is the opened mouth, (d) is the opened eye, and (e) is the closed eye.

First, we cut out a 16×16 image from the original image to test if the mouth is recognized in the image by a human. The mouth was not identified clearly in the 16×16 image, because the 16×16 image is too small to recognize the mouth as shown in Fig. 2. Therefore, we cut out a 36×36 image from the original image again and tried to recognize the mouth. With a 36×36 image, the mouth was identified clearly in the image. Lastly, we cut out a 52×52 image from the original image again and tried to recognize the mouth. Although we can also recognize the mouth in the 52×52 image, a 52×52 image contains redundant information such as noise. We cannot represent the features of mouth simply in our system by 52×52 images because of the noise.

As a result, we concluded that a 36×36 image is the best image to recognize the mouth (see Fig. 2), and we cut out a 36×36 images as input data in our system. The coordinates of the mouth detected by BP is not always just on the center of the mouth because of the positional gap. Therefore, we cut out a 36×36 image scanning from a 52×52 image to verify the mouth independent of the positional gap between the mouth position in the image and the coordinates of the mouth detected by BP. Nine 36×36 images are cut out from the 52×52 image scanning it by each 8 pixels from end to end, and twenty-four 36×36 images are cut out by scanning the 52×52 image by each 2 pixels around the center of the 52×52 image as shown in Fig. 3.

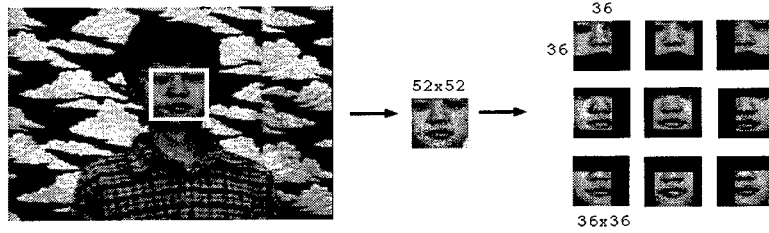


Fig. 3. A method of sampling a mouth.

In our system, the ratios of brightness contrast R_N ($N = 1 \sim 9$) are calculated for each area around the mouth based on Eq.1. The method for calculating the ratio is explained using Fig. 4. These ratios are fed to MAM described in the next section. The ratios are calculated in the different size images of mouths due to different distances from the camera to a person and individual variation. We calculated the ratio changing the size and position in the 36×36 image three times. In the result, we calculated six ratios, $(R_1, R_2) \times 3$, representing the right corner of the mouth, six ratios, $(R_3, R_4) \times 3$, representing the left corner of the mouth, and fifteen ratios, $(R_5, R_6, R_7, R_8, R_9) \times 3$, representing the nose. These values (R_1, \dots, R_9) are assigned to three-layers as the input vectors as shown in Fig. 5 where BR_N defines the brightest pixel in part-N, and DA_n defines the darkest pixel in part-N in Eq.1

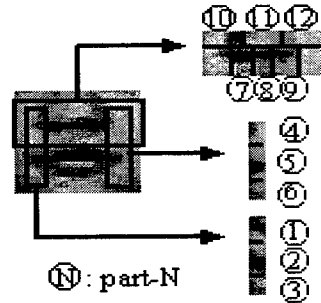


Fig. 4. Template for calculating the ratio.

$$\begin{array}{lll}
 R_1 = DA_2 \div BR_1 & R_2 = DA_2 \div BR_3 & R_3 = DA_5 \div BR_4 \\
 R_4 = DA_5 \div BR_6 & R_5 = DA_9 \div BR_8 & R_6 = DA_7 \div BR_8 \\
 R_7 = DA_{10} \div BR_8 & R_8 = DA_{12} \div BR_8 & R_9 = BR_{11} \div BR_8
 \end{array}$$

1

VERIFYING THE EYE USING MAM

MAM [9] is one of the associative memory neural networks, which recalls a pattern by linking several restored patterns together (see Fig. 5.). First, four training patterns $(X_1, Y_1, Z_1), \dots, (X_4, Y_4, Z_4)$ were prepared as shown in Table 1. These patterns are the bipolar data groups. The three patterns (pattern-1,2,3) represent the features of a mouth, and pattern-4 represents the features of non-mouth parts. The values of three training patterns (pattern-1, 2, 3 in Table 1) are set as bipolar, that is, 1 or -1. If -1 is assigned as the value of pattern-4, some training patterns cannot be recalled since there is correlation between the four training patterns (pattern-1~4). Therefore, we assigned -1.9 as the values of pattern-4, and pattern-4 is recalled [10].

Next, the sets of three vectors (A, B, C) were prepared as input data to MAM: Vector A represents the right corner of the mouth and vector B represents the left corner while vector C represents the nose. Then one of the training patterns was recalled. In our simulation, we assigned bipolar data to a 6-D vector A , a 6-D vector B , and a 15-D vector C (see Fig. 5.). The values of vectors, $(a_1, \dots, a_6, b_1, \dots, b_6, \text{ and } c_1, \dots, c_{15})$, were calculated based on the condition of R_n in Table 3. We assigned 1 to the values of vectors, (A, B, C) , if R_n is within the range shown in Table 3, otherwise we assigned -1 to the values of vectors, (A, B, C) . In the same way, bipolar data was assigned to values of three vectors, (A, B, C) , as in Fig. 5. In our system, the properties of brightness around the mouth are represented by these vectors.

Table 1. Training patterns for three mouth images of different sizes and a non-mouth image; Vector X is stored in layer-1, vector Y is stored in layer-2, vector Z is stored in layer-3 of MAM.

Pattern	Vector														
Pattern for large mouth image 1	X_1	1	1	-1	-1	-1	-1								
	Y_1	-1	-1	-1	-1	1	1								
	Z_1	1	1	1	1	1	-1	-1	-1	-1	-1	-1	-1	-1	-1
Pattern for medium mouth image 2	X_2	-1	-1	1	1	-1	-1								
	Y_2	-1	-1	1	1	-1	-1								
	Z_2	-1	-1	-1	-1	-1	1	1	1	1	1	-1	-1	-1	-1
Pattern for small mouth image 3	X_3	-1	-1	-1	-1	1	1								
	Y_3	1	1	-1	-1	-1	-1								
	Z_3	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	1	1	1	1
Pattern for non-mouth image 4	X_4	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9								
	Y_4	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9								
	Z_4	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9	-1.9

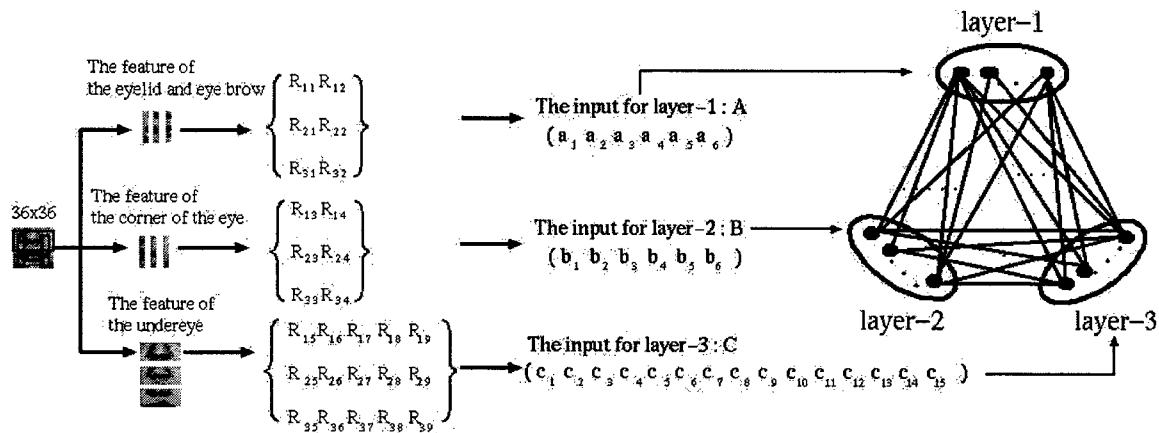


Fig. 5. Input vectors for three layer multidirectional associative memory.

Table 2. The range of R_n ($R_1 \sim R_9$ were calculated based on Eq.1.)

a_n	The range of R_n	b_n	The range of R_n	c_n	The range of R_n
a_1	$0.30 \leq R_1 < 0.55$	b_1	$0.30 \leq R_3 < 0.55$	c_1	$0.30 \leq R_5 < 0.60$
a_2	$0.30 \leq R_2 < 0.55$	b_2	$0.30 \leq R_4 < 0.55$	c_2	$0.30 \leq R_6 < 0.60$
				c_3	$0.60 \leq R_7 < 0.85$
				c_4	$0.60 \leq R_8 < 0.85$
				c_5	$0.95 \leq R_9 < 1.25$





SIMULATION RESULTS

In the simulation, we used three images of mouth and an image of eye extracted by BP as shown in Table 3. When (A, B, C) was given to MAM as an initial state, MAM converged to the steady state (A_f, B_f, C_f) . Image-1 and Image-2 represent different sizes of the mouth. Image-1 and Image-3 represent the mouths in different lighting conditions. The brightness of lighting is 500 lux for Image-1, and 58 lux for Image-3. We also tested the mouth image in dark lighting condition (18 lux), and our system can hardly distinguish

between mouth and eye because of the weak contrast of brightness of the mouth and nose. When the brightness of lighting is 18 lux in the image, the rate of correct recall decreased to 12.5%.

We prepared five patterns, (templates 1 to 5 as shown in Table 5. as templates of mouths to verify whether the input images include the mouth or not. If the steady state coincides with one of the five mouth templates, we determine the input image as a mouth. Template-1, 2, and 3 represent the training patterns. Template-4 and Template-5 represent the patterns, one of which is recalled when the input image includes the properties of two different mouth sizes. In all three cases of the mouth (image-1 ~3), one of the mouth templates is recalled by MAM as the steady state. In the case of the eye, (image-4), non-mouth template pattern-4 is recalled by MAM as the steady state. It took about only 1 second for this simulation. The simulation program was written in C language and simulated on SUN Ultra2 ($2 \times 168\text{MHz}$). We prepared 40 eye patterns and 40 mouth patterns as test data in all simulations, and the rate of correct recall is 813%.

Table 3. Simulation Results; Vectors (A, B, C) are the input for each layer of MAM, and vectors (A_f, B_f, C_f) are the steady states of MAM.

Image	State	Vector														
	Initial	A	-1	-1	1	1	-1	-1								
		B	1	1	1	1	-1	-1								
		C	-1	-1	1	1	-1	1	1	1		1	-1	-1	1	-1
	Steady	A_f	-1	-1	1	1	1	1								
		B_f	1	1	1	1	-1	-1								
		C_f	-1	-1	-1	-1	-1	1	1	1	1	1	1	1	1	1
	Initial	A	1	1	-1	-1	-1	-1								
		B	-1	-1	1	1	1	1								
		C	1	1	-1	1	-1	1	-1	-1	-1	1	-1	-1	-1	-1
	Steady	A_f	1	1	1	1	-1	-1								
		B_f	-1	-1	1	1	1	1								
		C_f	1	1	1	1	1	1	1	1	1	1	-1	-1	-1	-1
	Initial	A	-1	-1	1	1	-1	-1								
		B	-1	1	1	1	-1	-1								
		C	1	-1	-1	1	1	1	1	1	-1	-1	-1	-1	-1	-1
	Steady	A_f	-1	-1	1	1	-1	-1								
		B_f	-1	-1	1	1	-1	-1								
		C_f	-1	-1	-1	-1	-1	1	1	1	1	1	-1	-1	-1	-1
	Initial	A	-1	-1	-1	-1	1	1								
		B	1	1	-1	-1	-1	-1								
		C	-1	-1	-1	-1	-1	1	1	-1	-1	-1	-1	-1	-1	-1
	Steady	A_f	-1	-1	-1	-1	-1	-1								
		B_f	-1	-1	-1	-1	-1	-1								
		C_f	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1

DISCUSSION

It is very difficult to identify a mouth and eye completely in a small image, because the small image does not contain a conclusive clue to identify the mouth and eye. We cannot distinguish the mouth from the eye in 16×16 image, since the difference between eye and mouth is not obvious as shown in Fig. 2. Particularly, the closed eye tends to be mistaken for the mouth in the images. Therefore, the small image is inadequate to design the recognition system of the facial parts.

Table 4. Five templates of the mouth.

template	layer	Vector															
Template-1	1	1	1	-1	-1	-1	-1										
	2	-1	-1	-1	-1	1	1										
	3	1	1	1	1	1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
Template-2	1	-1	-1	1	1	-1	-1										
	2	-1	-1	1	1	-1	-1										
	3	-1	-1	-1	-1	-1	1	1	1	1	1	-1	-1	-1	-1	-1	-1
Template-3	1	-1	-1	-1	-1	1	1										
	2	1	1	-1	-1	-1	-1										
	3	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	1	1	1	1	1	1
Template-4	1	1	1	1	1	-1	-1										
	2	-1	-1	1	1	1	1										
	3	1	1	1	1	1	1	1	1	1	1	-1	-1	-1	-1	-1	-1
Template-5	1	-1	-1	1	1	1	1										
	2	1	1	1	1	-1	-1										
	3	-1	-1	-1	-1	-1	1	1	1	1	1	1	1	1	1	1	1

By analyzing the various sizes of images including a mouth, we found that 36×36 is the best image for the proposed system. First, the system detects the contrast of brightness on the periphery of the mouth. The contrast of brightness on the periphery of the visual field is perceived at the primary visual area in our brain. Since the contrast of brightness is discriminated in a digital image with low resolution regardless of different lighting conditions, it becomes an important clue to recognizing the mouth and eye in our system. After detecting the contrast of brightness, the system extracts the features such as both corners of the mouth, and the nose. Then, MAM recalls the mouth by linking these features based on the spatial context of the features, which are represented by the contrast of brightness[11].

Recent studies of visual recognition reveal that a human can subjectively recognize an object by linking related features on the periphery of the object when he or she cannot clearly identify it due to noise. Visual recognition by linking related features has been a controversial issue among many neural network researchers, and various models have been proposed to date. However, most have not been applied for practical use because the behavior of the models using chaotic or oscillatory dynamics is too sensitive to control completely, and they all take prohibitive computational time due to their complexity. In our system, the features around mouth are simply represented by the contrast of brightness, and the features are combined with each other by MAM to recall the mouth. The advantage of MAM is to recall a syntax vector in short computational time.

Finally, we address an interesting simulation result. When we expanded the size of image from 36×36 to 52×52, we found that the performance of the system worsened. We surmise that the 52×52 image contains redundant information as noise which is not useful to the network (see Fig. 2). A suitable size of the visual field depends on the value of information around the object, which is called “*the effective visual field*”. Note that it should be determined carefully to find the best size of effective visual field.

CONCLUSION

By extending the visual field for recognizing the mouth and linking the related features around the mouth, the associative memory model MAM can verify whether the target image includes a mouth or not, which is extracted by BP. The system performs well for the different color of lips, the lighting condition, and the state of the mouth, opened or closed mouth, by using the contrast of brightness around the both corner of the mouth, the tip of nose, and the nostril. A whole system became more complete by verifying the mouth in addition to the eye.

REFERENCES

1. Susumu Ohkita, et. al., 1998. "Detection Of Human Facial Parts using Infrared and Visible Images". EANN98, Gibraltar, Jun. 54-57.
2. Souichi Oka, et. al., 1998. "A self-Organized Oscillation in Detecting Time-Varying Human Faces", EANN98, Gibraltar, Jun. 58-61.
3. M. Ikeda and S. Saida, 1978. "Span of recognition in reading". Vision Res., 18, 83-88, 1978.
4. Mitsuo Ikeda, 1982. "A Recent Topic in Visual Psychophysics-Pattern Recognition and Visual Field Size". IEICE, vol65, 1288-1291.
5. Keith Rayner, 1978. "Eye Movements in Reading and Information Processing". Psychological Bulletin, Vol. 85, N0. 3, 618-660.
6. S. Ullman, G. Schechtman, 1982. "Adaptation and gain normalization". Proceedings of the Royal Society of London, B 216, 299-313.
7. David Marr, 1982. "Vision: A computational Investigation into the Human Representation and Processing of Visual Information". San Francisco, CA: Freeman, Chapter3.
8. Miki Kitabata, et. al., 1998. "Confirmation of the Detected Eye by Bidirectional Associative Memory". EANN98, Gibraltar, Jun. 70-73.
9. M.Hagiwara, 1990. "Multidirectional Associative Memory". IJCNN-90-WASH-DC, I, 3-6.
10. Wang Y.F., Cruz, Jr. J.B., Mulligan Jr. J.H., 1990. "Two Coding Strategies for Bidirectional Associative Memory". IEEE Trans. on Neural Networks, 1, 81-92.
11. Irving Biederman, 1972. "Perceiving Real-World Scenes". Science, 177, 77-80.

A Current-Mode Sorting Circuit for Pattern Recognition

Gu Lin and Bingxue Shi

Institute of Microelectronics, Tsinghua University, Beijing, 100084, P.R.China

ABSTRACT

A current-mode sorting circuit based on magnitude for pattern recognition is proposed. In this sorting circuit, symmetric WTA (Winner Take All) network is employed to find maximum current. Then, sorted currents based on magnitude are outputted in time-shared way. This sorting circuit has a simple and flexible structure. Results of experiment show that it has good performances. It can be widely used in pattern recognition, classification, expert system and so on.

INTRODUCTION

Pattern recognition is one of the important parts in artificial intelligent field. Generally, pattern recognition process simply includes two steps[1,2,3]. The first one is that unknown pattern is matched with each exemplar pattern to calculate the matching scores. The matching scores are a measure of similarity between the input pattern and the exemplar patterns. The second one is that the exemplar pattern with the maximum matching score is selected. That is, the exemplar pattern, which is nearest to the input unknown pattern, is outputted as recognition result. Thus, in these pattern recognition hardware system based on above process steps, only one exemplar pattern with the maximum matching score is selected as recognition result. However, with greatly increasing the complexity of system and the number of standard patterns, especially, development of multistage cascade systems, the way selecting only one exemplar pattern could not meet requirements of system performances. So it is necessary to find out two or more exemplar patterns nearer to the unknown pattern according to the matching scores, this requires that the system is able to sort the matching scores based on their magnitudes, instead of maximizing the matching scores.

At present, many sorting integrated circuits have been proposed, however, almost of them are digital sorting circuit[4,5,6,7]. Although digital sorting circuits can also used for analog signals, A/D and D/A converters are required, so that the structure of circuits will be even more complicated. In addition, errors may occur in the conversion between digital signals and analog signals. To this end, a novel analog current-mode sorting circuit that is able to be easily implemented in VLSI is proposed in this paper.

The proposed sorting circuit is able to output the input currents in a sorting order on the basis of the magnitudes of the input currents. In addition, the sorting circuit can find out the order of input terminals based on the magnitudes of the input currents. The sorting time is in linear relation to the number of input currents. With simple structure, the sorting circuit can be widely used in pattern recognition[8,9], classification, expert system and so on.

STRUCTURE AND OPERATION PRINCIPLE OF THE SORTING CIRCUIT

For convenience, a current-mode sorting circuit based on magnitude is discussed by taking three input currents for example. The circuit diagram of the sorting circuit is shown in Fig.1. The circuit diagram of TRANS unit in Fig.1 is shown in Fig.2. This sorting circuit comprises following four blocks.

*This project(69636030) is supported by National Natural Science Foundation of China

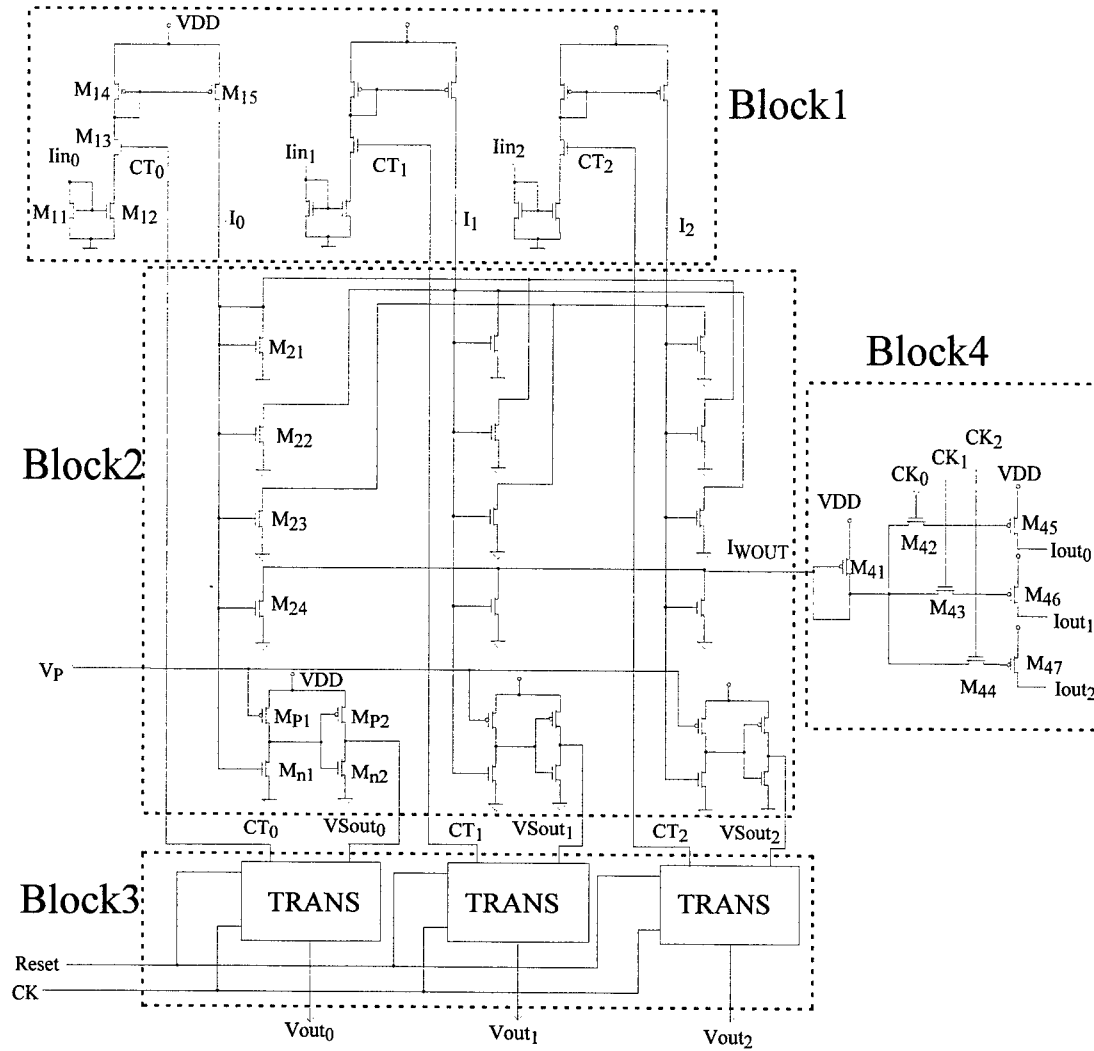


Fig.1 The circuit diagram of the current-mode sorter with three input currents

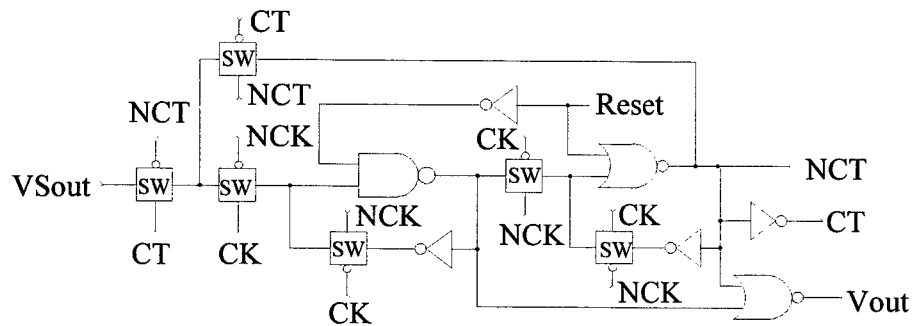


Fig.2 The circuit diagram of the feedback control and voltage output unit

Block1 is an input block composed of three fully identical input units. $I_{in0} \sim I_{in2}$ designate three input currents to-be-sorted. For the sake of convenience, only the input unit on the left-most side is described. In this input unit, NMOS current mirror with unit gain is composed of NMOS transistors M_{11} and M_{12} , and PMOS current mirror with unit gain is composed of PMOS transistors M_{14} and M_{15} . A switch transistor M_{13} controls the magnitude of an output current I_0 . When the switch transistor M_{13} is on, I_0 equals I_{in0} , and when the switch transistor M_{13} is off, I_0 is zero.

Block2 is current-mode WTA (Winner Take All) circuit network. This WTA circuit is a laterally inhibitory interconnected network with high resolution and high speed. Assuming $I_0 = \max(I_0, I_1, I_2)$. After inputting currents $I_0 \sim I_2$, the network executes convergence operations and finally reaches a steady status. The maximum input current is outputted as I_{wout} , namely, $I_{wout} = I_0 = \max(I_0, I_1, I_2)$. At the same time, a high level of output voltage is obtained on node $VSout_0$, and low voltage level is outputted on node $VSout_1$ and $VSout_2$. Namely, the function of finding maximum input current is completed.

Block3 is a feedback control and voltage output circuit which is composed of three fully identical circuit TRANS units. The circuit diagram of the TRANS unit is shown in Fig.2, where the SW unit therein is a CMOS switch and the NCK is an inverse signal of CK. The feedback control and voltage output circuit generates feedback control signals $CT_k (0 \leq k \leq 2)$ according to $VSout_k (0 \leq k \leq 2)$, which are outputted from WTA circuit, to control the output currents of the input block1. In addition, the feedback control and voltage output circuit is able to convert the low-to-high voltage level from $VSout$ to a high output voltage pulse from $Vout$. This high voltage pulse is used to determine the corresponding input terminal with respect to the sorted output current for processing the sorted currents.

Block4 is the output block for outputting sorted currents. In the block, switch transistors M_{42}, M_{43} and M_{44} are respectively controlled by non-overlapped clock signals CK_0, CK_1 and CK_2 . PMOS transistors M_{45}, M_{46} and M_{47} are mirror transistors, each is identical to M_{41} in size. Under the control of clock signals $CK_k (0 \leq k \leq 2)$, the output current I_{wout} from the WTA circuit can be mirror-mapped to output terminals in a time shared way to generate I_{out_0}, I_{out_1} and I_{out_2} . The currents I_{out_0}, I_{out_1} and I_{out_2} are the sorted currents of input currents $I_{in_k} (0 \leq k \leq 2)$.

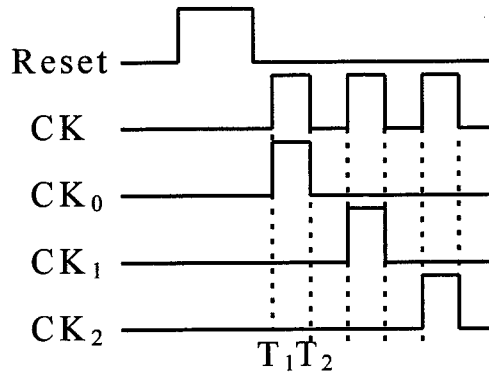


Fig.3 The timing diagram of sorting circuit

The operation principle of the sorting circuit with three inputs is considered as follows. Fig. 3 shows timings of all clock signals. Firstly, the reset signal Reset increases from low level to high level, this causes the level at the node $CT_k (k=0,1,2)$ in block3 goes to high, and the level at the node $Vout_k (k=0,1,2)$ goes to low. Since the level at the node $CT_k (k=0,1,2)$ is high, $I_k (k=0,1,2) = I_{in_k} (k=0,1,2)$ in block1. Then, the WTA circuit works. For the convenience of description, herein I_{in_0} is assumed as the maximum input current, that is, $I_{in_0} = \max(I_{in_0}, I_{in_1}, I_{in_2})$. Therefore, the WTA circuit finds the maximum input current by lateral inhibitory effect and outputs the maximum input current at I_{wout} , that is, $I_{wout} = I_{in_0}$. At the same time, this causes the level at node $VSout_0$ to be high and the levels at node $VSout_1$ and $VSout_2$ to be low. At the instant of T_1 , since the clock signal CK_0 goes to high, the transistor M_{42} in block4 is turned on, so that the maximum current I_{wout} is mirror-mapped to the drain of M_{45} to generate I_{out_0} , that is, $I_{out_0} = I_{wout} = I_{in_0}$. When the clock signal CK_0 goes to low, I_{out_0} is still the maximum current I_{in_0} due to the sample/hold effect of the switched current mirror. On the other hand, at the instant of T_1 , the clock signal CK goes to high. It is clear that the $Vout_0$ becomes high since the level at the node $VSout_0$ is high, while $Vout_1$ and $Vout_2$ become low since the levels at the node $VSout_1$ and $VSout_2$ are low. When the clock signal CK goes to low at the instant of T_2 , the level at the node CT_0 becomes low, and $Vout_0$ becomes low, so that a high voltage pulse is generated on the $Vout_0$ terminal. While CT_1 and CT_2 are still high, $Vout_1$ and $Vout_2$ are still low. Among $Vout_0, Vout_1$ and $Vout_2$, only $Vout_0$ outputs one high level pulse, this also shows the terminal of

input current I_{in0} has the maximum input current value. Since the level at the node CT_0 is low, the voltage at the node $VSout_0$ can not be inputted to the block3 so that the level at the node CT_0 and $Vout_0$ will hold zero until the next reset signal is activated. On the other hand, the levels at the node CT_0 being low makes the current I_0 become zero in block1, so that I_0 will not influence the sequential comparing operation. In this manner, I_{in1} and I_{in2} are compared and the maximum one between them will be determined by the process described above and held on $Iout_1$ terminal. A high voltage pulse is also generated on the corresponding $Vout$ terminal. The remaining operations can be deduces accordingly. Finally, the sorted results can be obtained.

From above discussion, in block4, all of the input currents to-be-sorted are presented on $Iout_k(k=0,1,2)$ in an order of magnitudes under the control of the clock signals. Meanwhile, high voltage pulses are sequentially generated on the corresponding $Vout$ terminals for determining the input terminals with respect to the sorted currents. In addition, this sorting circuit is flexible and the output way could be easily controlled by adjusting the combination ways of clock signal $CK_i(i=0,1,2)$. For example, when only CK_0 is inputted, the sorting circuit is simplified as MAX circuit. When only CK_2 is input, the sorting circuit is simplified as MIN circuit.

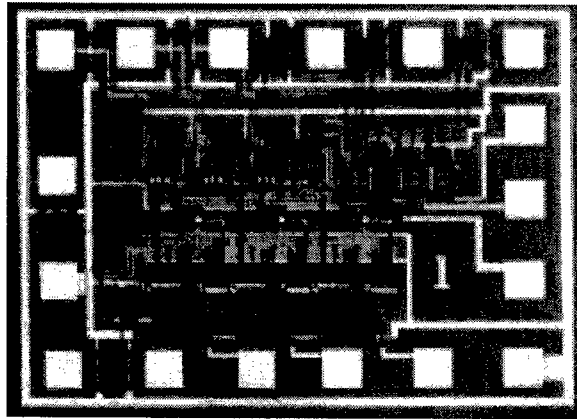


Fig.4 The microphotograph of the sorting circuit

EXPERIMENT RESULTS

The sorting circuit has been successfully fabricated in $2\mu\text{m}$ N-well standard digital CMOS process. The chip microphotography is shown in Fig.4. It has been successfully measured. The measured results are shown in Table.1 . Experimental results show the circuit has correct function and good performances.

Table 1. Measured Results for the current-mode Sorting circuit

frequency	higher than 5MHz
power supply	+ 5V
highest sorting resolution	$2\mu\text{A}$
average sorting resolution	$5\mu\text{A}$
sorting precision	better than $10\mu\text{A}$
range of input current	$10\mu\text{A}\sim 300\mu\text{A}$
power consumption	lower than 4mW

CONCLUSION

The proposed sorting circuit based on magnitude has good performance and high resolution. It has a simple, flexible structure. Because of switched-current structure employed, this sorting circuit is able to be directly fabricated in a standard digital CMOS process and be easily integrated in mixed analog-digital mode.

ACKNOWLEDGMENT

This project(69636030) is supported by National Natural Science Foundation of China

REFERENCES

1. U.Clingiroglu. 1993. A charge-based neural Hamming classifier. IEEE. J.Solid-State Circuit, 28(1),59-67.
2. Binqiao Li, Zhijian Li, Bingxue Shi, 1993. An analogue integrated circuit of a Hamming neural network designed and fabricated in CMOS technology, IJCNN93, Nagoya, Japan.
3. Gu Lin and Bingxue Shi , 1998. Novel multifunction switched-current fuzzy processor for pattern recognition, Chinese Journal of Semiconductors ,19(4), 291-297.
4. C.M.Blair. Low cost sorting circuit for VLSI. 1996. IEEE Trans on Circuits and Systems:Fundemental Theory and Application.43(6),515-516.
5. C.D.Thompson, 1983. The VLSI Complexity of Sorting. IEEE Trans on Computers. C-32(12),1171-83.
6. H.M.Alnuweiri, 1993. A new class of optimal bounded-degree VLSI sorting networks. IEEE Trans on Computers. 42(6),746-751.
7. A.R.Seigel, 1985. Minimum Storage Sorting Circuits. IEEE Trans.Computers,C-34(4),355-361.
8. Gu Lin and Bingxue Shi, 1999. A programmable and expandable hamming network integrated circuit, Proc.Second International Conference on Intelligent Processing and Manufacturing of Materials(IPMM99) ,Honolulu, Hawaii.
9. Gu Lin and Bingxue Shi, 1999. A multi-input current-mode fuzzy integrated circuit for pattern recognition, Proc. Second International Conference on Intelligent Processing and Manufacturing of Materials (IPMM99), Honolulu, Hawaii.

Intelligence in the Design of Materials and Processes II

Intelligent Design Methods for Smart Materials

M. Fathi-Torbaghan, L. Hildebrand

University of Dortmund, Department of Computer Science, Chair 1
Otto-Hahn-Str. 16, 44227 Dortmund 50, Germany

ABSTRACT

The design process of modern smart materials often require the use of complex system models. These models cannot be derived easily due to the complex knowledge that describe the process. In some cases, model parameters can be gained using neural networks, but these systems allow only a one-way simulation from input values to learned output values. If evaluation in the other direction is needed, these models allow no direct evaluation. This task can be solved using modern techniques like evolutionary algorithms and fuzzy logic. The use of such a combination allows evaluation of the learned simulation models in the direction from output to the input. An example can be given from the field of screw rotor design.

INTRODUCTION

In recent years fuzzy logic has become widely acknowledged – apart from its other applications – as an important and useful methodology in the design of rulebased systems [1, 2, 3]. It allows the representation of imprecise or incomplete knowledge and offers various mechanisms for reasoning with fuzzy data. In comparison to “classical” rulebased systems, only very few rules are needed to describe difficult problems. Nevertheless, in its current form it has several shortcomings: when it comes to the design of membership functions or to actually attaching priorities to the available rules, the choice of numerical quantities for the different parameters which is indispensable for the reasoning process is generally not justified by the results from knowledge acquisition and, what is worse, demands often a long process of iterative improvement to obtain good results. The use of empirically obtained quantitative representations seems questionable because of its high context dependence. The results are in many cases sub-optimal systems.

It seems natural to try to use a computer and an algorithmic optimization technique for the final adjustment of the parameters. To our mind, evolutionary algorithms seem especially appropriate for this task, partly because the fuzzy reasoning process can hardly be described by means of a closed mathematical formula – not to mention differentiability or other “convenient” mathematical properties –, partly because of the opportunity to apply parallel computation in a very natural way which seems essential in the design of large-scale systems. The combined use of artificial intelligence methods, like expert systems and blackboard architectures, and computational intelligence methods, like evolutionary algorithms and fuzzy systems, lead to a powerful approach for complex industrial tasks.

CONSTRUCTION OF SMART MATERIALS

Fuzzy logic can be used to model the knowledge of construction of composite materials. But before this, we would like to explain the idea of composite materials. In the last few years, new technologies have influenced the requirements of industrial production. New materials are needed to fulfill these demands. A mismatch between requirements and material properties can cause drastic consequences, with a cutback of function or resorting to superior materials with properties that are overqualified for the requirements. The solution is the design of composite materials. Only these materials guarantee economic fulfillment of requirements in order to be competitive in the industry. Examples for the use of composite materials are:

- engine parts manufactured using ceramics to decrease the weight of moving components,
- automotive industry, e. g. corrosion-resistant coating of steel as corrosion prevention,
- tool making, e.g., very firm tool surfaces to increase sharpness and hardness of tools such as drills or chisels.

To produce these materials, different techniques such as embedding particles, fibers, whiskers and coating or laminating are used. These three different types of composite materials are shown in figure 1:

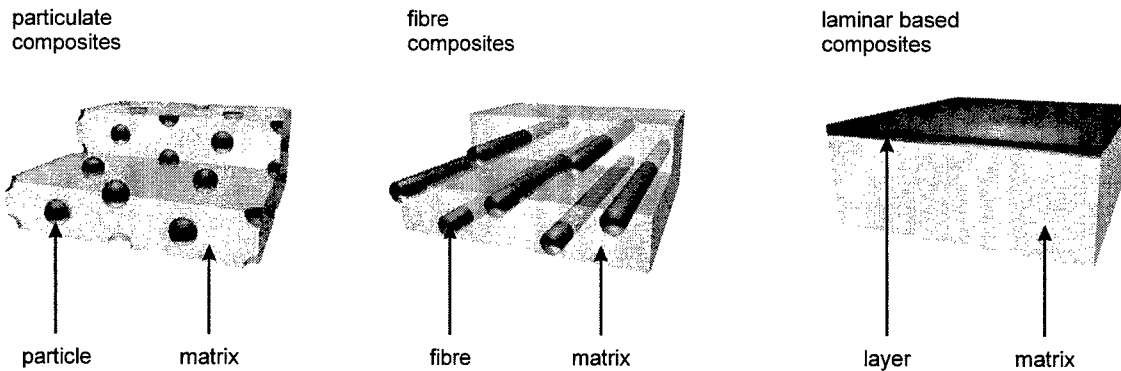


Fig. 1. Structure of composite materials

The proper selection of specific components for composite materials requires adequate tools for the predetermination of material properties. Especially the design of screw rotor is a hard problem. Screw rotors are the most important mechanical part of various types of screw rotor machines. These machines can be used to produce compressed air, as well as a special kind of motors, the screw motors. Figure 2 shows an example of a screw compressor, where the different screw rotor pairs can be seen. Many attempts have been done to increase the performance and reliability of the surface of the screw rotors. One important aspect is the covering of ceramics to improve the wear behavior of the whole system [4, 5].

The design of the screw rotor profiles is a very complex task. Both rotors need to be mutually inverse, the distance of the rotating parts must be in accurate specified limits and a lot of mechanical properties must be fulfilled. Figure 3 gives an overview of the design process. The first step is to determine a mathematical model from which the screw rotor profiles are generated. If this is satisfactory, the rotor are build into the machine and the whole system has to be tested. Depending on the type of system, a compressor or a motor, different condition must hold.

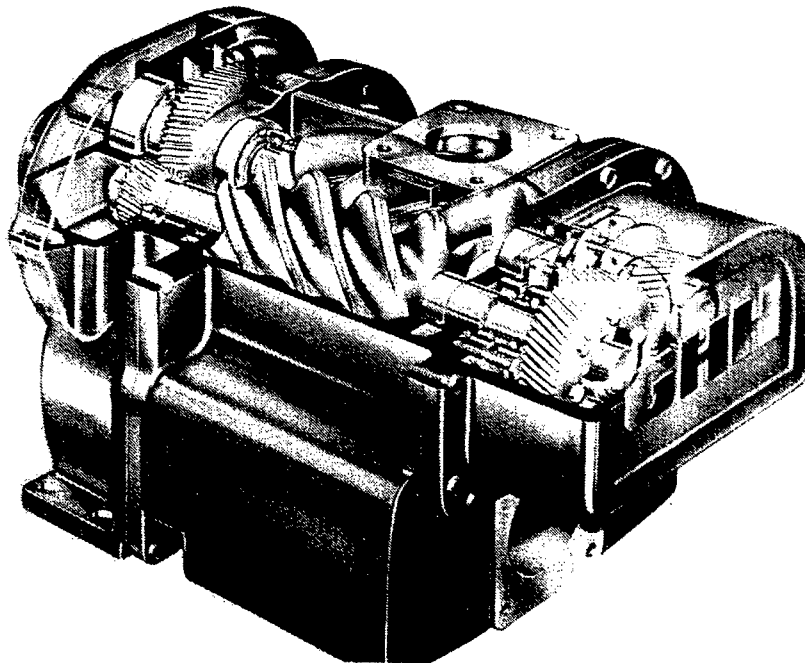


Fig. 2. Screw compressor with build in screw rotors.

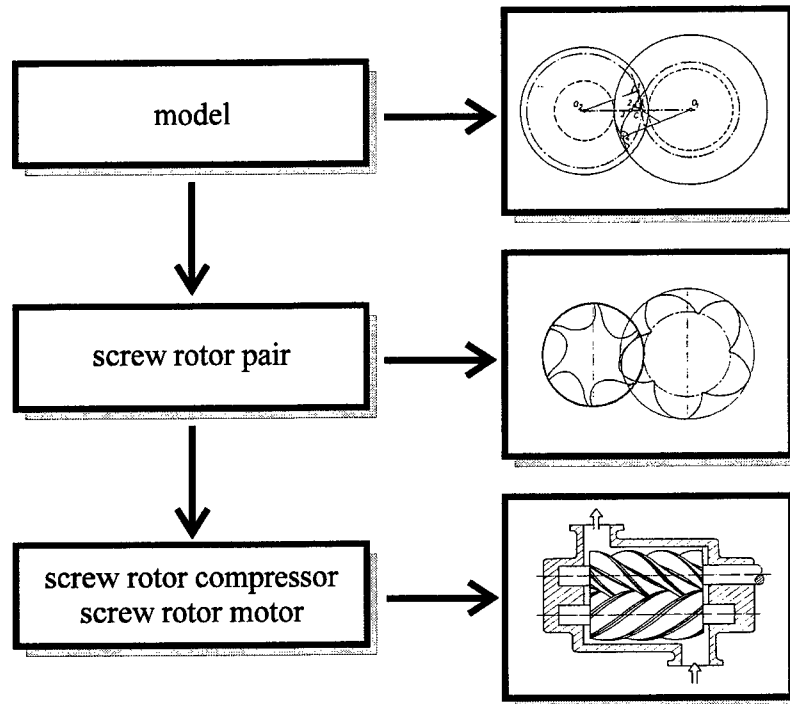


Fig. 3. Design steps of a screw rotor.

Numerous theoretical and empirical models have led to unsatisfactory results because applicability of these models is very restricted. A new way to evaluate more universal models is to use optimizing methods from the research field of computational intelligence [4, 5, 6, 7].

INTELLIGENT METHODS

KEE was used as the primary problem solving method in the blackboard based system EXCOCOM (EXpert system for COstruction of COmposite Materials). The kernel of EXCOCOM consists of several blackboards, each partitioned into several levels of abstractions, and of a set of knowledge sources. Some of these knowledge sources are linked to external systems, i.e. to a FEM (Finite Element Methods) package or to a database. The functions of a blackboard are varied. They establish the communication and interaction between the knowledge sources and enable a global repository for incremental generation of the solution. EXCOCOM uses six blackboards which are now described in detail:

- blackboard 1 handles the user input. This input is checked for completeness and missing data is calculated (if possible) or demanded from the user.
- blackboard 2 decides which technology (homogeneous or composite) should be selected.
- blackboard 3 searches for possible composite materials
- blackboard 4 and 5 support this decision making
- blackboard 6 simulates the composite materials found with the FEM-module

To gain control over the different blackboards and knowledge sources a special control system for the information exchange between the blackboards had to be implemented. Due to the fact that the expert system KEE is a problem solving system based on a theoretical logical model we had problems modeling the experts knowledge. Mathematical or physical knowledge such as formulas or facts could easily be implemented, but all kinds of uncertain and inexact knowledge such as heuristics and experiences could not be used using the blackboard approach.

Conventional expert systems have a theoretically oriented starting point to handle expert behavior. They are based on Boolean logic, so they cannot describe reality, which is uncertain and indistinct. The results of the knowledge acquisition of designing composite materials were general rules with linguistic terms, and many

facts and data are missing. In the second approach, we will show that fuzzy expert systems are able to represent these indefinite concepts. To implement the fuzzy expert system the main concepts of fuzzy logic for rulebased systems have been used. Only small changes have been made to start the development with a well-defined system. To enter knowledge into the system in the form of membership functions and fuzzy rules, a fuzzy editor and a rule editor are put at the user's disposal. The inference engine is capable of working with backward and forward chaining so that the inference engine can activate or deactivate used or unused knowledge. On the basis of this architecture, the fuzzy expert system will be able to solve different tasks automatically if these tasks are related to each other. The inputs to the system are the limits which the materials must withstand and the system designs the composite material synthetically.

Evolutionary algorithms form a class of probabilistic optimization techniques motivated by the observation of biological systems [8, 9, 10, 11]. Although these algorithms are only crude simplifications of real biological processes, they have proved to be very robust and due to their simple parallelizability even efficiently implementable. The basic idea of evolutionary algorithms is the use of a finite population of individuals, any one of which represents exactly one point in search space. After its initialization the population develops in a self organizing collective learning process constituted by the (stochastic) subprocesses of selection, mutation and recombination towards better and better regions of the search space. The selection process favors those individuals with a high fitness value, the mutation operator supplies the population with new genetic information, and the recombination operator controls the mixing of this mutated information when it is passed from one generation to another. The use of evolution strategies is necessary because the experts are not able to formulate all rules required and to give an exact description of the membership functions. Due to these inaccuracies a method to optimize the fuzzy rule based system is useful. The decision to apply evolution strategies has been made to improve model fuzzy membership functions and to fill in gaps in the rule base [5, 7].

In many cases knowledge acquisition is more like a translation of knowledge into rules and membership functions than a process of gaining new knowledge. For this reason it differs from the classical knowledge acquisition. At first glance it seems to be an easy task to do this translation manually. But past experience has shown that experts as well as specialists of fuzzy logic are not able to determine the rules and membership functions in the required accuracy. Modern optimization methods from the field of computational intelligence seem to be a suitable basis for the task of function approximation [7, 12, 13]. Its main part is an optimization loop, which assesses the current fuzzy system and tries to improve it using an evolution strategy until the difference between the input data and the output data of the fuzzy system is less enough. Input is the exact knowledge in form of tables and diagrams. An initial fuzzy system is generated and evaluated over the entire range of definition. The next step adds up the differences between the input data and the output data of this fuzzy system to calculate a single value, which describes the accuracy of the current system. The evolution strategy modifies this system in the next step to find an improved, which means a more accurate, one. This optimization loop continues until a sufficient overall accuracy is obtained.

CONCLUSION

The process of gaining knowledge can be reduced to entering input data and starting an automated process. Function approximation using fuzzy logic and evolutionary strategies is able to acquire exact knowledge without the help of experts. This feature makes it possible to acquire the huge amount of knowledge which already exists in form of tables and diagrams. This kind of representation is commonly used in the field of designing screw rotors. Combined with classical acquisition methods for vague and imprecise knowledge the introduced methodology provides a good basis for the implementation of knowledge from other technical oriented fields.

Early approaches to model the knowledge developed in the special investigation area have shown many disadvantages. They were too slow, not very user-friendly or not accepted by the experts. Since the introduction of the fuzzy system we have increased the acceptance of a computer based system as well as the speed of implementing new knowledge. To our mind a fuzzy-system is a good approach to model the knowledge which has been found during the last years. The amount of knowledge and the easy

implementation make it clear, that a fuzzy system is able to model knowledge of small subproblems as well as the knowledge of large scale projects.

REFERENCES

1. L.A. Zadeh, 1971. Quantitative fuzzy semantics. Inform. Sci., 3.
2. L.A. Zadeh, 1979. A theory of approximate reasoning. in L.I. Mikulich J.E. Hayes, D. Mitchie, eds., Machine Intelligence, 9, 149-194. Wiley, New York.
3. A. Kandel, Ed., 1992. Fuzzy Expert Systems. CRC Press.
4. M. Fathi, 1993. A fuzzy expert system for real-time monitoring in manufacturing, Proc. 12th World congress international federation of automatic control - IFAC.
5. M. Fathi, L. Hildebrand, 1999. Complex System Analysis using CI-Methods, in: Proceedings "AeroSense - SPIE 13th Inter. Symp. on Areospace/Defense Sensing, Simulation, and Controls", Orlando, USA.
6. D. Dubois, H. Prade, R. R. Yager, 1993. Fuzzy sets for intelligent systems, Morgan Kaufmann, San Mateo.
7. M. Fathi, L. Hildebrand, 1997. "Model-free optimization of fuzzy rulebased systems using evolution strategies", IEEE Trans. on System, Man, and Cybernetics, Part B: Cybernetics, 27(2).
8. J. Rechenberg, 1973. Evolutionsstrategie: Optimierung technischer Systeme nach Prinzipien der biologischen Evolution. Frommann-Holzboog Verlag, Stuttgart.
9. H.-P. Schwefel, 1994. Evolution and Optimum Seeking, John-Wiley & Sons, New York.
10. J. H. Holland, 1975. Adaption in Natural and Artificial systems. Ann Arbor, MI. U. Michigan Press.
11. D. E. Goldberg, 1989. Genetic Algorithms in search, optimization and machine learning. Reading, MA: Addison-Wesley.
12. M. Fathi, L. Hildebrand, 1994. Evolution Strategies for the Optimization of Fuzzy Rules. Conf. IPMU, Paris.
13. M. Fathi, L. Hildebrand, 1995. The application of evolution strategies to the problem of parameter optimization in fuzzy rulebased system. IEEE International Conf. on Evolutionary Computing. Perth.
14. L. Hildebrand; M. Jäger, M. Fathi, 1998. Learning of Linguistic and Numerical Knowledge - Application of Neural Networks and Evolutionary Algorithms, in: Proceedings "Engineering of Intelligent Systems EIS98", La Laguna, Spain.
15. M. Fathi, L. Hildebrand, 1998. Using Uncertainty for the Design of Composites, in: Proceedings "Third Pacific Rim International Conference on Advanced Materials and Processing" (PRICM-3), Honolulu, HI.
16. L. Hildebrand; M. Fathi, 1998. Evolutionary Design of Screw Rotor Profiles, in: Proceedings "International Conference on Composite Engineering" ICCE/5 '98, Las Vegas, USA.

Identification of a Model which Relates Variations in Shape Geometry to Discrete Process Control Variables in Shape Forging

B.F. Rolfe^{*}, M.J. Cardew-Hall^{*}, S.M. Abdallah^{*}, G.A.W. West^{}**

^{*} Dept. of Engineering, The Australian National University, Canberra, Australia, 0200

^{**} School of Computer Science, Curtin University of Technology, Perth, Australia, 6102

ABSTRACT

This paper develops a model that identifies changes in process control parameter values by analysing the geometry of a product from a forging process. That is, shape variations from the norm of a forged part are linked to the variations of the process parameters of the forging process. This is an important aspect of inspection analysis because it implies that the analysis of shape variation can return the potential errors within the product and also the changes needed to the process control parameters to correct errors in the product shape. Our model can identify whether input parameters of interest are too high or too low to an accuracy of at least 70%. Given the complexity of the forging process these initial results show promise.

INTRODUCTION

There is inherent variability in any manufacturing system and the forging process is no different in this respect. What we as engineers hope to achieve is to control and reduce the variability we see in the process. To do this we must first characterise what variability there is and second find what parameters assist in reducing this variability. This paper develops a model that identifies changes in process control parameter values by analysing the geometry of a product from a forging process. This investigation follows on from Daniel's *et al.*[3] work which examined geometric variations in forging to aid in design and tolerancing of die sets.

Our model has been developed using the following steps:

- Create a set of forged parts;
- Analyse the set of parts to find the main modes of variation;
- Learn the relationship between the modes of variation and the process control parameter values.

The forging process has been simulated using finite element analysis (FEA) software. The process control parameter values of the forging process were varied about a mean set of values to create a set of forged billet shapes. These were analysed using a deformable model based on principle component analysis (PCA). This analysis was used to characterise the variations within the shape into major variation modes. Each shape in the set of forging parts was characterised by a weighting vector which combines varying amounts of each major variation mode. The set of weighting vectors were used to create a response surface with respect to the varying process control parameter values. A classifier was used to determine the levels of the process control parameters from the response surface. In an abstract sense the classifier was used to learn the shape of the response surface to identify the original process control parameter values. To simplify the classification we divided each process control parameter's response into three regions: high; normal; and low.

It should be noted, however, that this analysis of the final geometry assumes that varying the different input parameters of the process will produce distinct and different end geometries, otherwise the analysis will be intractable without some further knowledge. Similarly, this model relies on the inductive process, that is, it learns from specific examples and then abstracts to the general. The specific examples must therefore give enough information about the general system otherwise the model will fail.

DESCRIPTION OF THE FORGING PROCESS

A finite element analysis (FEA) model was used to simulate the forging process. This made changing the input parameters easier when creating a set of forged shapes for the training set. Some of the input parameters to the FEA model were varied to affect changes in the end shape of the billet. Figure 1 shows the geometry of the forging process. Note that material data was not varied as an input into the forging process.

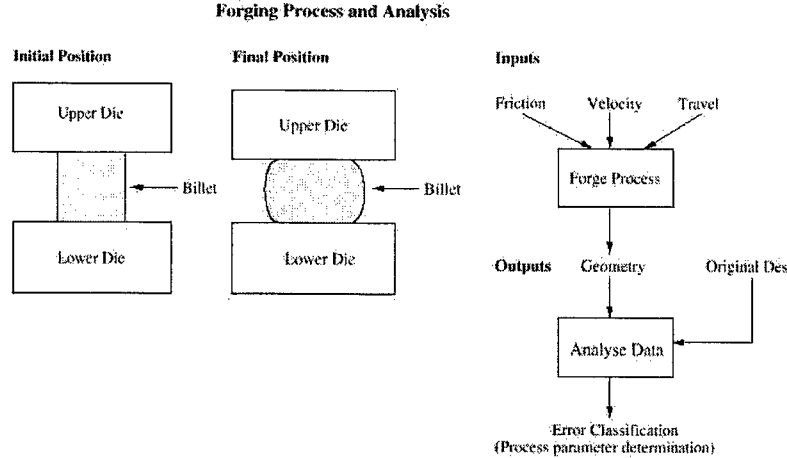


Fig. 1. The forging process and analysis schematic.

A medium carbon steel (0.45% C) was selected for the billet material because it is often used in hot working. In general the hot-working flow stress equation is a power function of strain rate $\dot{\epsilon}$ [9]. The simulation package used for the simulations was Forge 2.7¹. This package calculates each iteration by relating the deviatoric stress tensor, σ , to the strain rate tensor, $\dot{\epsilon}$ as follows:

$$\sigma = 2 K (T, \bar{\epsilon}) (\sqrt{3} \dot{\epsilon})^{m-1} \dot{\epsilon} \quad 1.$$

where T is the temperature in degrees Kelvin and $K(T, \bar{\epsilon})$ is the strength or consistency equation.

There is a lot of flexibility in the ability to vary many parameters within the forging simulation package, but due to the complexity of the problem at hand we have limited our scope to only three of the main forging variables: friction between the punch and the billet; velocity of the punch; travel of the punch. We were able to see changes in the cross sectional geometry of the forged billet by varying the input parameters of the process. Initially a set of "perfect" parameter values was chosen that would produce the "perfect" shape. We chose the following "perfect" parameter values: friction = 2.0×10^{-1} ; velocity = 10.0 mm/sec; travel = 20.0 mm. A set of input parameters was created by varying each of the parameters up to $\pm 50\%$ of their initial set value. Because of the sheer number of simulations involved, the set was limited to one and two parameter variations for any particular simulation.

SHAPE ALIGNMENT

After creating a set of training shapes, each shape in the set needed to be made consistent. Each shape is just a list of boundary nodes which describe the boundary surface of the deformed billet. All shapes created for this paper were roughly aligned before forging and maintained their alignment throughout processing. The boundary nodes of each shape were made consistent by the following process: determining a constant datum for every shape; increasing or decreasing the number of boundary nodes to maintain the same number of boundary nodes for all shapes. A consistent datum point was determined across all shapes. This datum was then used as the start point for each list of boundary nodes for every shape. The list of boundary nodes was also made consistent across all the shapes by ensuring each list had the same number of boundary nodes (200 nodes).

¹ Forge 2.7 is a forging simulation package owned by Transvalor S.A.

PRINCIPLE COMPONENT ANALYSIS

Once all the shapes were aligned a principle component analysis (PCA) was implemented on the set of shapes using the coordinates of the boundary nodes as data. PCA realises a restricted set of eigenvectors which describe most of the variations or deformations in the set of shapes.

PCA Background

For this paper we use a type of statistical deformable model, *point distribution model* (PDM), developed by Cootes *et al.* [2] based on PCA. They used PCA to describe modes of variation in two dimensional data when examining heart images[1]. Later, the model was extended to three dimensional data to also analyse heart data[5, 6]. Daniel *et al.* [3] used the *point distribution model* to inspect both 2D and 3D forging data.

PCA is a statistical technique that finds the directions of maximum variability inherent in the data set. It is based on eigen structure analysis of variations from a training set of shapes. The resulting eigenvectors correspond to the major modes of variation. The eigenvectors are found by creating a covariance matrix from the training set of shapes. The covariance matrix is given by,

$$S = \frac{1}{N-1} \sum_{i=1}^N (\hat{X}_i - \bar{X})(\hat{X}_i - \bar{X})^T \quad 2.$$

where $\hat{X}_i = [x_{i1}, y_{i1}, z_{i1}, \dots, x_{ki}, y_{ki}, z_{ki}]^T$ is the i^{th} shape coordinates and is the mean shape of all the shapes in the training set. Moreover, (x_{ij}, y_{ij}, z_{ij}) is the i^{th} boundary point on the j^{th} shape in the training set. Finding the eigenvectors of S realises the principle components of the training set. The eigenvectors are sorted in descending order by their corresponding eigenvalues. Thus, the most significant principle components are the first few eigenvectors. By rule of thumb the most significant principle components which explain 90% of the variations are chosen and placed into matrix P . So, any shape in the training set can be well represented by

$$\hat{X}_i = \bar{X} + P b \quad 3.$$

where b is the weighting vector showing how much of each principle component is needed to vary the mean shape to the shape \hat{X}_i .

Forging PCA

After all the shapes are aligned, PCA is performed on the set of shapes. We update the PCA equation (3) for the forging process as follows:

$$\hat{X}_{pca} = \hat{X}_{mean} + P b \quad 4.$$

where \hat{X}_{pca} is the list of boundary point coordinates for any shape and \hat{X}_{mean} is the list of boundary point coordinates for the mean shape of the training set. The matrix P is found from calculating the statistical variations within the shapes of the training set (see equation (2)). The principle component analysis realises a restricted set of eigenvectors that describe most of the variations in the set of shapes. From these components we determine the corresponding weighting vectors, b , which deform the *mean shape* into every shape in the training set. Rearranging equation (4) we get,

$$b = P^+ (\hat{X}_{measured\ shape} - \hat{X}_{mean}) \quad 5.$$

where the P^+ is the pseudo inverse of P , and \hat{X}_{mean} is the mean shape also defined in the identification phase. The variable $\hat{X}_{measured\ shape}$ contains the points on the surface of the newly produced part. The resulting b vector is then analysed by the parameter identification model defined in the section on identifying process parameters.

Find the major modes of deformation

PCA relies on finding the major modes of deformation from statistical data obtained from the training data set. The "perfect" shape defined previously was used as the *mean shape* as defined in the forging PCA equation (4). For each of the six data sets a covariance matrix S was calculated. The P matrix was created from the six most significant deformation modes, that is, the six most significant eigenvectors. The first six modes were chosen after examining the results from using two, four, six and eight modes. After choosing

the number of modes, the b vectors for every shape in every data set were calculated using the pseudo inverse of the P matrix using equation (5).

RESULTS OF PCA AND THE ASSOCIATED RESPONSE SURFACES

Component Plots

The first two components of the calculated b vectors are plotted vs. each other in Fig. 2 for two data sets.

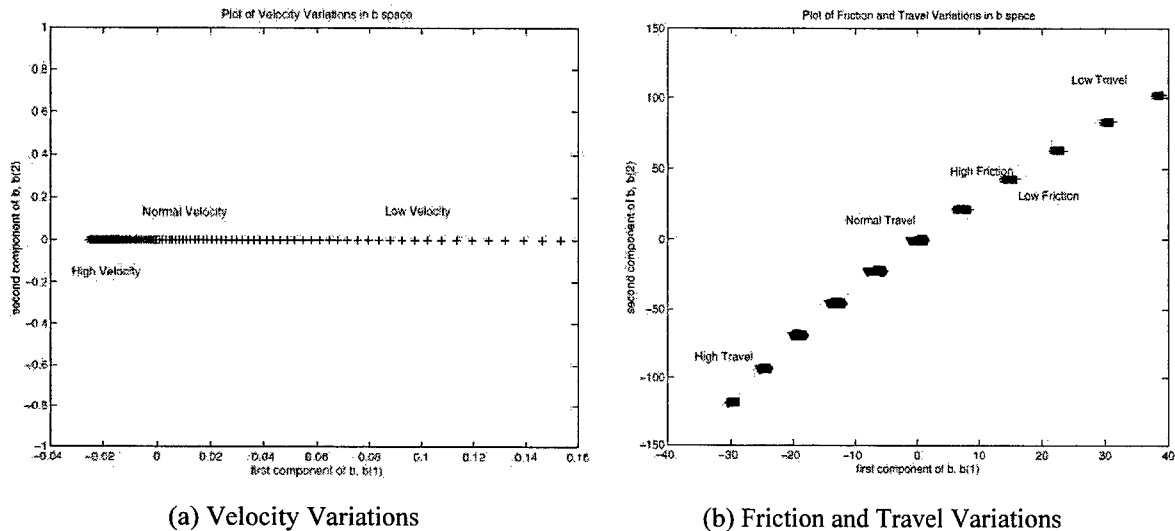


Fig. 2. Parameter Variations in b space.

The single line of points in the velocity single variation data set indicates that velocity variations create a set of single weighted versions of a particular shape deformation, i.e., there is only one mode of variation. This was also seen in the other single variation data sets (friction and travel). The velocity data set has an interesting phenomenon of the high velocity data points being squashed together while the low velocity data points are spread out. This means the lower the velocity of the punch, the greater the change on the end shape according to this particular set-up for the simulated forging process. The friction/travel data set shows that the travel parameter is dominant. The shape points move from the top right corner to the lower left corner as the travel value is increased, whereas increasing friction just moves the shape points from right to left around the given travel value. Note that each parameter variation data set has its own particular deformation mode and therefore each b space cannot be directly compared with another.

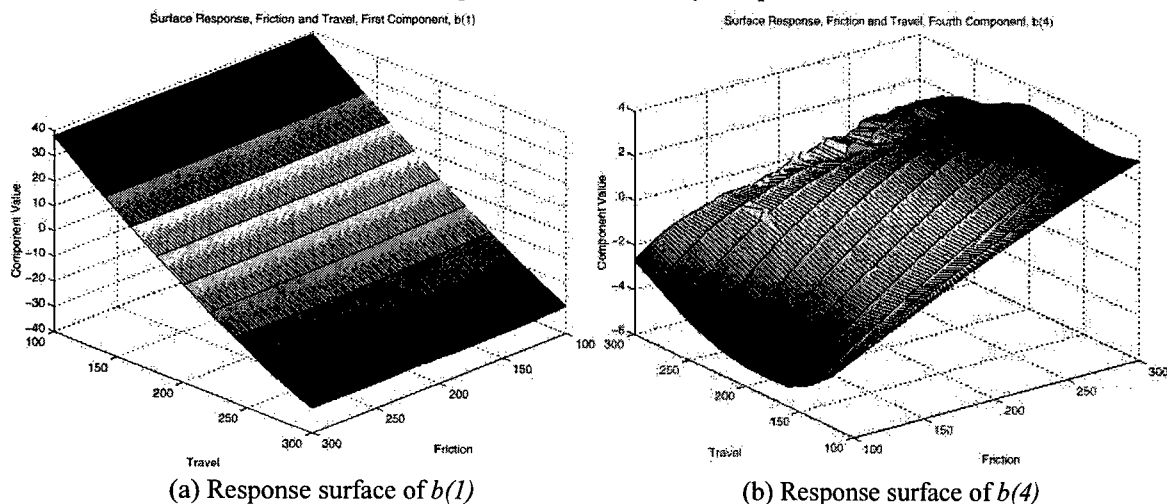


Fig. 3. Response surfaces of b with respect to friction and travel.

Response Surfaces

The data sets containing the calculated b vectors gave a response surface with respect to the input parameters. We plotted the surface for a single component of the b vector with respect to several varying input parameters. This gave a visual indication of the variation in the b component with respect to the varying input parameter. For example, Figure 3(a) is the response surface of the most significant b component, $b(1)$, with respect to friction and travel. As we have noted previously, the travel effect is much greater than the friction effect. This can be seen by the fact that the surface ranges from -40 to +40 when travel is varied, however, the surface ranges approximately from -1 to +1 when friction is varied. On the other hand, because friction is a complex phenomenon it has a much greater effect at higher orders (less significant) b components. Figure 3(b) shows how friction has a more dominant effect than the travel. The surface ranges from -6 to +4 when friction is varied and approximately -6 to -3 when travel is varied.

Dividing the input parameter space into regions

It was initially felt that it would be easier to classify b space if each parameter was separated into a series of regions. We split the input parameter space into three separate regions for each process parameter (high, normal, low). This gave three regions for one parameter and nine regions for two parameters. The boundaries were chosen arbitrarily and without bias towards the data. Splitting the regions in this fashion allowed us to have a qualitative output which, in the first instance, was easy to interpret. Furthermore, then coarseness of the regions could be reduced in the future to provide more quantitative information about the input parameters being observed.

IDENTIFYING PROCESS PARAMETERS

Primarily we were interested in two things: first, could we identify variations in one parameter; second, could we segregate multiple variations in more than one parameter. When we split input parameter space into regions what we are doing is segmenting b space into regions. Each region in input parameter space has a corresponding region in b space. Unfortunately, the components of the b space do not directly relate to any known form of product error, but rather, to a combination of known product errors. Classification techniques were then used to discriminate between the different levels of the process variables.

Learn to discriminate between regions

There are several types of discriminator that could be used, statistical, rule based, decision trees or artificial neural networks (ANN)[8]. In our case we did not have much statistical information from the FEA simulation. This was because a given set of input parameters to the forging process would always produce the same shape and therefore the same b vector. Thus, the ANN was eventually chosen over the other three methods because of its ease of use and its ability to segregate non-linear data. Initially we tried a simple two node single output layer ANN which was enough to segregate single or dominant parameters into low, normal or high regions, such as travel. However, the interesting problem was to classify data that had two or more varying parameters. Due to the non-linearities involved when varying two parameters the simple neural network was not able to converge to a solution.

Classifying the data using the trained discriminator

A single hidden layer feedforward backpropagation ANN was chosen to classify the data [8]. After a quick sensitivity analysis it was determined that the single hidden layer should have double the number of neurons (18) as the number of outputs (9) where the number of output neurons equals the number of regions.

To determine how well we can train the ANN with respect to each data set, we conducted a series of 10-fold cross over tests. The 10-fold cross over test consists of choosing 100 random sample points from the data set. Initially the first 10 are retained for testing and the remainder are used for training. This is repeated for a second 10 samples while the remaining 90 samples are used for training and so on until there have been 10 separate and distinct tests. The results of the 10-fold validation are combined into a "confusion matrix" [4]. The "confusion matrix" is a $N \times N$ contingency table of actual region versus classified region where N equals the total number of regions. If the ANN was a perfect classifier then the confusion matrix would be a diagonal matrix. The value of the diagonal represent misclassification. The confusion matrix is used to

indicate the accuracy of the classifier with respect to available data. The accuracy to classify a randomly chosen sample was determined by dividing the number of correct classifications by the total number of points classified for each column of the confusion matrix. Low accuracy implies difficult classification. The results in Table 1 were generated by running several 10-fold tests and taking the average of the results.

Table 1: Classification Results.

Region	Travel Friction	Velocity Friction	Velocity Travel
	% Accuracy	% Accuracy	% Accuracy
Low Param 1 Low Param 2	72.7	90.4	80.5
Med Param 1 Low Param 2	77.8	77.9	72.7
High Param 1 Low Param 2	75.0	90.7	79.6
Low Param 1 Med Param 2	84.5	73.6	57.6
Med Param 1 Med Param 2	76.6	68.0	69.2
High Param 1 Med Param 2	80.8	80.3	65.9
Low Param 1 High Param 2	85.2	71.9	76.6
Med Param 1 High Param 2	70.5	79.5	84.6
High Param 1 High Param 2	69.3	77.9	79.1

One way to see how well the classification works is to compare the accuracy of the classifier to a "random assignment" classifier [7]. Random assignment has an accuracy of one in every nine samples, that is, 11.1%. The Friction Travel data set gives good results with correct classifications occurring 70 to 80 % of the time. The Friction Velocity data set gave slightly better results with accuracies for each region above 77%. The Travel Velocity data set gave reasonable accuracies. The regions of medium travel appear to be difficult to classify correctly with accuracies of 57.6%, 69.2% and 65.9%. Probably this is due to the points within these regions being very close and so the classifier cannot easily discriminate the differences in velocity. This is supported by the spread of misclassifications for these regions in the confusion matrix. This situation can be changed if more information can be found to discriminate between the different velocity regions.

CONCLUSIONS

We have introduced a method to characterise variations in shape by using a point distribution deformable model. PDM uses principle component analysis to determine major modes of variation. These major modes then give a quantifiable way to measure shape. Furthermore, by training an ANN, we have created a model which can determine, to an accuracy of at least 70%, whether input parameters are too high or too low. This model can be used for inspection purposes to determine quantitatively, how incorrect a part may be. Also, identification of whether input process control variables are too high or low can contribute to a control system that determines what measures can be taken to correct a situation. Our method to develop a model to identify input process control variables from shape needs to be tested on real forge data to complete this investigation. However, the initial results from finite element analysis simulations show great promise.

ACKNOWLEDGMENTS

The authors would like to thank the support of an Australian Research Council grant number A49532454 and would also like to acknowledge that ANU has a major project with Ford Australia, STAMP, which is investigating ways to improve efficiency and quality of sheet metal processed car body parts.

REFERENCES

1. T. F. Cootes, et al., 1992. A trainable method of parametric shape description. *Image Vision Comp.*, 10, 289-294.
2. T. F. Cootes, et al., 1995. Active shape models - their training and application. *Computer Vision and Image Understanding*, 61(1), 38-59.
3. B. T. Daniel, et al., 1997. Geometric Variations: Analysis, Optimisation and Control, in *Spatial Computing: Issues in vision, multimedia and visualisation techniques*. World Scientific: Singapore, 79-114.
4. G. H. Dunteman, 1984. *Introduction to Multivariate Analysis*. Beverley Hills: Sage Publications.
5. A. Hill, et al. 1992. A generic system for image interpretation using flexible templates, *British Machine Vision Conf.*, Springer-Verlag.
6. A. Hill, et al. 1993. Model-based interpretation of 3D medical images, *British Machine Vision Conf.* BMVA Press.
7. M. James, 1985. *Classification Algorithms*, Collins, London.
8. R. Schalkoff, 1992. *Pattern Recognition: statistical, structural and neural approaches*, J. Wiley & Sons, Toronto.
9. J. A. Schey, 1987. *Introduction to Manufacturing Processes*. New York: McGraw-Hill Book Company.

Mechanical Characteristics of Hipped SiC Particulate-Reinforced Aluminum Alloy Metal Matrix Composites

C.Y. Chung and K.C. Lau

Department of Physics and Materials Science,
City University of Hong Kong,
Tat Chee Avenue, Hong Kong

ABSTRACT

Al-Cu alloy matrix composites reinforced with various volume fraction of SiC particulates (SiCp) were prepared by conventional powder metallurgy (PM) and hot isostatic pressing (HIP) processes. The tensile and fracture behavior of PM and HIPed composites were studied. The experimental results were compared with the Shear-lag, Eshelby and modified Eshelby micro-mechanics models. The composite stiffness tends to increase with increasing SiCp volume fraction. The modified Eshelby model correlated well with the experimental results whereas the stiffness prediction of shear-lag model was lower than the experimental data. The fracture stresses of PM and HIPed metal matrix composites decrease with increasing SiCp volume fraction. This was attribute to the cracking of SiCp upon tensile deformation.

INTRODUCTION

Metal matrix composites (MMCs) have attracted considerable attention in recent years as good structural materials. These hybrid composites inherent the metallic (high ductility and toughness) and ceramic (high strength and modulus) characteristics and exhibit unique advantages over conventional alloys. Aluminum MMCs are attractive because of their high specific modulus and strength, superior thermal stability, and improved wear resistance. MMCs are now being increasingly used in components for aerospace, automotive and power industries [2,3].

MMCs are usually reinforced by ceramic phase, which may be in the form of fibers, whiskers or particles. The continuous long fiber reinforced MMCs are expensive with strong an-isotropic properties. Discontinuously reinforcement usually give lower strength and lower stiffness, but they are isotropic and compatible to the conventional metal-forming processes. The properties of MMCs depend on the volume fraction, morphology, size-distribution of the reinforcements, microstructure of the matrix, and the interfacial properties between the ceramic and matrix phases. The properties of the matrix dominate when the reinforcement volume fraction is low. The casting and powder metallurgy (PM) techniques generally give discontinuously reinforced MMCs. These processes are of lower cost, good surface finishing and high reproducibility. However, machining the cast MMCs is difficult owing to the abrasive resistance nature of the ceramic reinforcement. Furthermore, the non-wetting of ceramic-metal interface, incomplete infiltration, and segregation of reinforcements are common problems for the MMC casting process. In these aspects, the PM process giving near-net-shape products is more advantageous for the fabrication of whisker or particulate reinforced MMCs. The PM fabricated MMCs are usually of small sub-grain size microstructure with limited segregation of particles. These contribute to the improved mechanical properties of many PM fabricated materials [1-6].

Cavities or pores very often exist in the sintered PM products. Additional mechanical treatments such as rolling, extrusion or hot pressing are needed to reduce the internal voids. Recently, hot isostatic pressing (HIP) technique is being increasingly used for the densification of PM products. The HIPping process consists of sintering the compacted powder mixture, which are sealed inside a capsule, at high pressure and high temperature. The pressurizing medium is normally inert gas such as argon. Several workers have successfully fabricated Al-based, Fe-based and Ti-based MMCs using the HIPping process [7-9]. The microstructure, tensile and failure characteristics of conventional PM fabricated Al-based MMCs are

well-documented [1-6,10-12]. Less information is available on the mechanical properties of the HIPped Al-based MMCs. Niklas et al. reported that the HIPped Al-based MMCs exhibit better mechanical properties relative to the extruded MMCs [7]. The aim of this paper is to study the properties of the HIPped and conventionally sintered Al-Cu MMCs with Silicon carbide particle (SiCp) reinforcement.

EXPERIMENTAL PROCEDURE

The composites system studied was Al-4wt%Cu matrix reinforced with SiCp. 36 μ m SiCp, 60 μ m pure aluminum powders and 50 μ m pure copper powders were used. The MMCs for this study were produced by conventional PM sintering and HIPping processes. Tensile bar specimens were produced at a compaction pressure of 350MPa after mixing and blending. The PM sintering was performed at 620°C for 1.5 hour in an argon atmosphere. The bars were further peak-aged at 150°C for 72 hours. For the HIPped MMCs, the powder mixtures were mixed and then cold compacted at 350MPa. HIPping was performed using an ABB mini-HIPper (Autoclave System QTH-3). The green compacts were consolidated at 600°C and 100MPa for 2 hours. The HIPped MMCs were also peak-aged at 150°C for 72 hours. Optical microscopy was used to observe the microstructure of MMCs. Immersion density measurements were carried out according to Archimedes' principle. In this technique density was determined by measuring the difference between the weight of a specimen in air and when it was submerged in distilled water at room temperature. Tensile measurements were performed using an INSTRON tensile tester (model 4206) at a crosshead speed of 1 mm/min.

RESULTS AND DISCUSSIONS

Fig. 1 shows the optical micrographs of the Al-4wt%Cu MMCs containing 20vol% of SiCp. SiCp distributed uniformly in the alloy matrix of the HIPped MMC, and there is no clustering or segregation of particles. On the other hand, numerous cavities are present in the matrix of the conventional sintered PM MMC. Fig. 2 shows the variation of the MMC density against the SiCp content. The densities of the HIPped MMCs are considerably higher than those of conventional PM sintered MMC. The densities of the HIPped MMCs are close to the fully-densified theoretical values up to 20vol% SiCp, and is lower than this theoretical density when volume fraction of SiCp \geq 25vol%. This implies that the HIPping treatment is effective in closing the pores, and thereby enhancing densification.

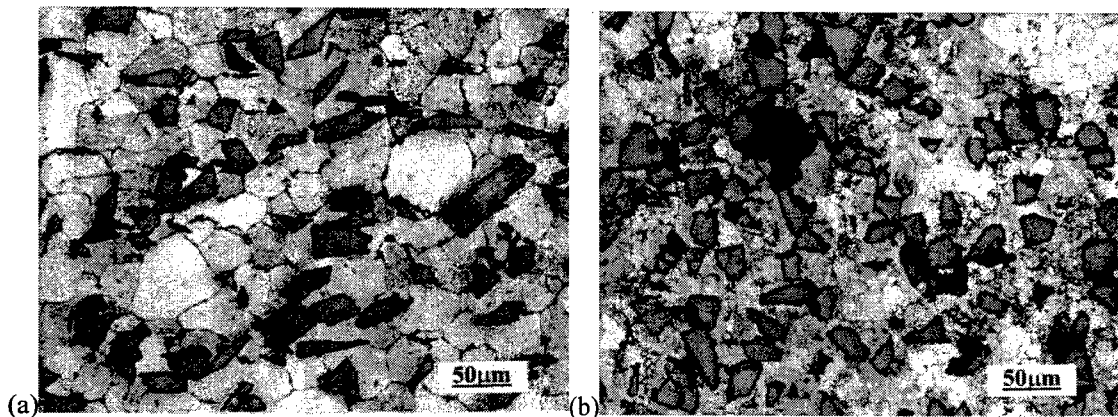


Fig. 1. Microstructure of MMCs containing 20 vol% SiCp fabricated by (a) HIPping and (b) conventional PM sintering processes.

Fig. 3 shows the stress-strain curves of the MMCs prepared by HIPping process. As expected, the yield strength of MMCs increase with increasing SiCp volume fraction whilst the elongation (ductility) decrease. The effect of discontinuous SiCp reinforcement on composite modulus can be predicted using micro-mechanics models such as Shear-lag and Eshelby. The Shear-lag model was originally developed by Cox [13], and generally used to calculate the stiffness and strength of short fibers or whiskers reinforced composites. This model assumes that all of the stress transfer from the matrix to fiber occurs

by interfacial shear around the periphery of the fiber. The modulus of the composites based on the Shear-lag model can be expressed as follows [14,15],

$$E_c = (1 - V_f)E_m + V_f E_f \left[1 + \frac{(E_m/E_f - 1) \tanh(\beta l/2)}{\beta l/2} \right] \quad 1.$$

where β is given by
$$\beta = \frac{2\sqrt{2}}{d} \sqrt{\frac{G_m/E_f}{\ln(D/d)}} \quad 2.$$

E_f , E_m and E_c are the Young's modulus of the fiber, matrix and composites, respectively; V_f is the reinforcement volume fraction; G_m is shear modulus of the matrix; D is width of unit cell for Shear-lag analysis; l and d are the length and the diameter of the fiber.

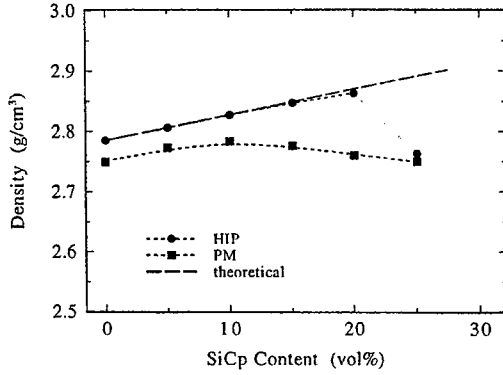


Fig. 2. Density vs. particulate volume fraction for the MMCs fabricated by HIPping and conventional PM sintering processes.

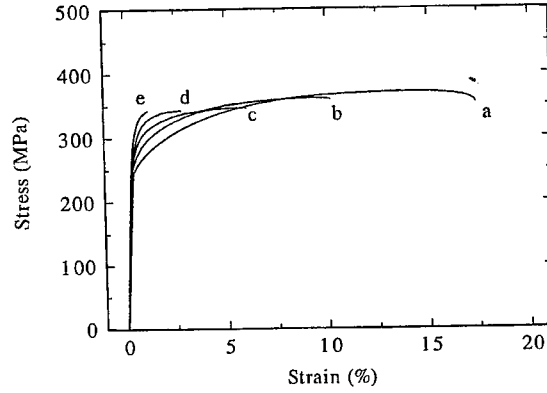


Fig. 3. Typical stress-strain curves of HIPped specimens (a) Al-4wt%Cu alloy, (b) MMC with 5%SiCp, (c) MMC with 10%SiCp, (d) MMC with 15%SiCp, (e) MMC with 20%SiCp.

For the SiCp reinforced MMCs, we assume that the particulates exhibit a spherical shape. Correspondingly, the aspect ratio (length to diameter) of the SiCp equals to unity,

$$V_f = d^3 / D^3$$

$$\beta = \frac{2\sqrt{2}}{d} \sqrt{\frac{3G_m/E_f}{\ln(V_f)}} \quad 3.$$

Fig. 4 shows the Young's modulus of reinforced MMCs against SiCp volume fraction. The stiffness of the composite determined from Shear-lag model (Equation (1)) is considerably lower than that of the experimental results. This is because the tensile load transfer at the fiber ends is ignored in the conventional Shear-lag model [16]. Furthermore, Shear-lag model tends to give a poor prediction of the stiffness of MMCs when the aspect ratio of the fiber reinforcement is smaller, or when the short fibers are mis-oriented. The aspect ratio of SiCp reinforcement is much smaller than that of the short fiber; hence, shear-lag model gives rise to lower modulus value, as shown in Fig. 4.

In the Eshelby model, composites consists of ellipsoidal inclusion (Ω) with non-elastic strain (ϵ_j^*) embedded in a continuous matrix. Such inclusion is stressed uniformly under general loading condition [17]. The ellipsoidal shape is a good approximation for real inclusion shape, i. e., an elongated ellipsoid simulates a fiber, and an equiaxed inclusion simulates a sphere. The stress inside Ω is given by

$$\sigma_j = C_{kl} (\epsilon_{kl} - \epsilon_{kl}^*) \quad 4.$$

where C_{ijkl} is the elastic constant tensor; ϵ_{kl} is the total strain and related to ϵ_j^* by

$$\epsilon_{kl} = S_{klmn} \epsilon_{mn}^* \quad 5.$$

where S_{klmn} is the Eshelby's tensor and it is a function of the geometry of the ellipsoidal inclusion and the Poisson's ratio (ν) of the matrix.

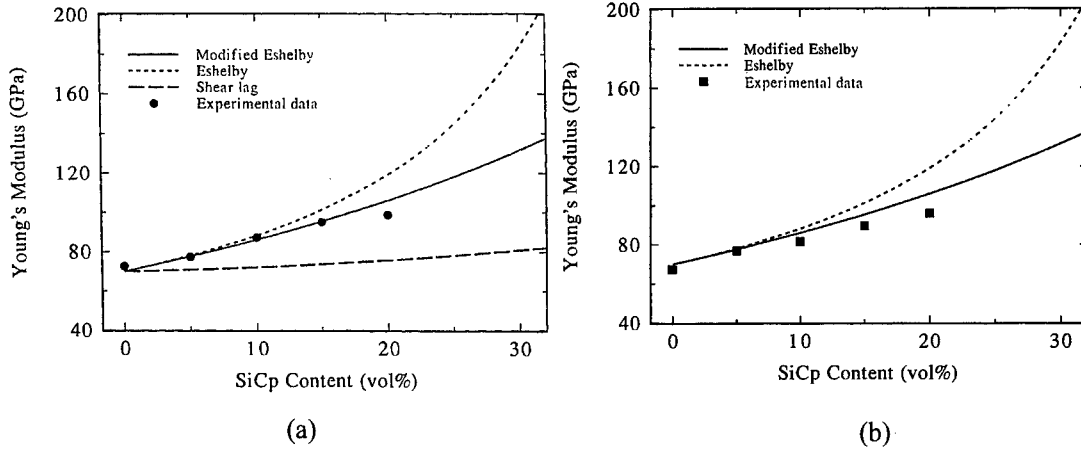


Fig. 4. Young's modulus vs. SiCp volume fraction for MMCs fabricated by (a) HIPPING and (b) conventional PM sintering process.

Young's modulus predicted from Shear-lag, Eshelby and modified Eshelby models are also shown.

For a spherical inclusion, Eshelby's tensor for elasticity [18] is given by

$$S_{1111} = S_{2222} = S_{3333} = \frac{7 - 5\nu}{15(1 - \nu)} \quad 6.$$

$$S_{1122} = S_{1133} = S_{2211} = S_{3311} = S_{3322} = \frac{5\nu - 1}{15(1 - \nu)} \quad 7.$$

$$S_{1212} = S_{2323} = S_{3131} = \frac{4 - 5\nu}{15(1 - \nu)} \quad 8.$$

For simplicity, we let $S_{1111} = S_{2222} = a$ 9.

$S_{1122} = S_{1133} = S_{2211} = S_{2233} = b$ 10.

From the above formulation, the stiffness of the spherical particulate reinforced composites [19] can be determined and expressed as:

$$E_c = \frac{2b^2 - a(a+b)}{2b^2 - a(a+b) + (a+b+2b\nu)V_f} \cdot E_m \quad 11.$$

The Young's modulus of the spherical particulate reinforced MMC estimated from the Eshelby model is also shown in Fig. 4. It is evident that the calculations from the Eshelby models are in better agreement with the experimental results for MMCs containing SiCp. Since Eshelby considered a simple ellipsoidal inclusion embedded in an infinitely elastic body, it is only be valid for small volume fraction reinforcement.

A more complex analysis is needed for high SiCp volume fraction, where interactions between the stress fields, associated with each inclusion must be taken into consideration. Mori and Tanaka [20] modified the Eshelby model to account for the integration between inclusions. Accordingly, the stiffness of the spherical particulate reinforced composite [19] can be expressed as,

$$E_c = \frac{(k+a-b)(k+a+2b)}{k(a-2b) + (a-b)(a+2b)} \cdot E_m \quad 12.$$

where $k = \frac{V_f}{1 - V_f}$ 13.

The modified Eshelby model gives better predictions of the Young's modulus for both SiCp reinforced MMCs fabricated by conventional PM sintering and HIPPING processes.

Fig. 5 shows the plots of 0.2% yield strength of MMCs against the SiCp volume fraction. The yield strength of the HIPPed MMCs are considerably higher than that of PM fabricated MMC because the HIPPING process has densified the MMCs due to the isostatic sintering condition. The yield strength increases with increasing SiCp volume fraction. The yield strength estimation of the composites using the Shear-lag model is also shown. The yield strength [14,15] can be expressed as

$$\frac{\sigma_c}{\sigma_m} = 0.5V_f(2 + l/d) + (1 - V_f) \quad 14.$$

where σ_c and σ_m are yield stresses of the matrix and composites, respectively. Similar to the Young's modulus calculations, the Shear-lag model give yield strengths lower than the experimental results.

The effect of SiCp reinforcement on composite fracture strength is shown in Fig. 6. The tensile strength decreases with increase in SiCp volume fraction. The SEM micrographs of the HIPped composites revealed that the MMCs fail through progressive particle cracking during tensile deformation when the volume fraction of SiCp is high, Fig. 7. This reveals that SiCp reinforcement has broken because the failure strain of particulates is smaller than that of the alloy matrix. The fracture stress of the composite can be estimated from the rule of mixtures, and assuming the fracture stress of particulate (σ_p) equals to zero, thus

$$\sigma_c = (1 - V_f) \sigma_{mf} \quad 15.$$

where σ_{mf} is the fracture stress of the matrix. The predicted yield stresses of MMCs fabricated by PM and HIPS techniques are depicted as full and dashed lines in Fig. 6. The experimental results correlate well with the theoretical estimation.

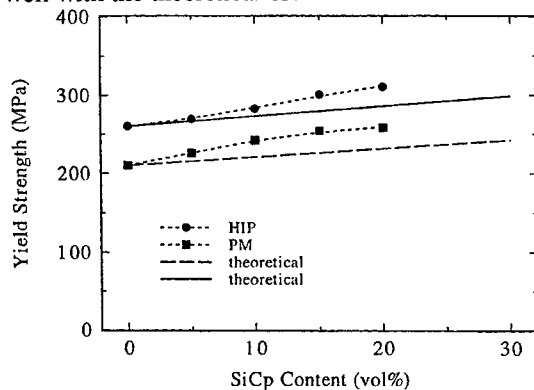


Fig. 5. Yield strength (σ_y) Vs SiCp volume fraction for HIPped and conventional PM sintering processes MMCs. The dashed and broken lines are yield stresses of the MMCs predicted from $\frac{\sigma_c}{\sigma_m} = 0.5V_f(2 + l/d) + (1 - f)$

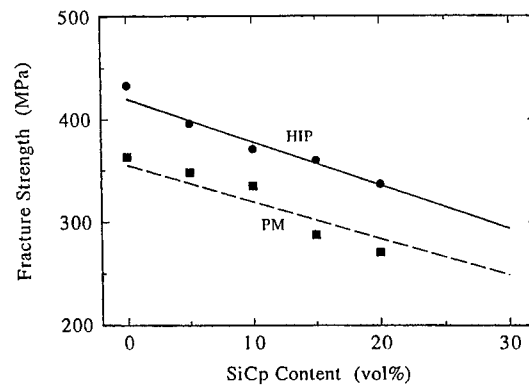


Fig. 6. Fracture strength of MMC vs. particulate volume content for MMCs fabricated by HIP and PM process. The dashed and broken lines and fracture stresses of MMCs predicted from equation $\sigma_c = (1 - f) \sigma_{mf}$

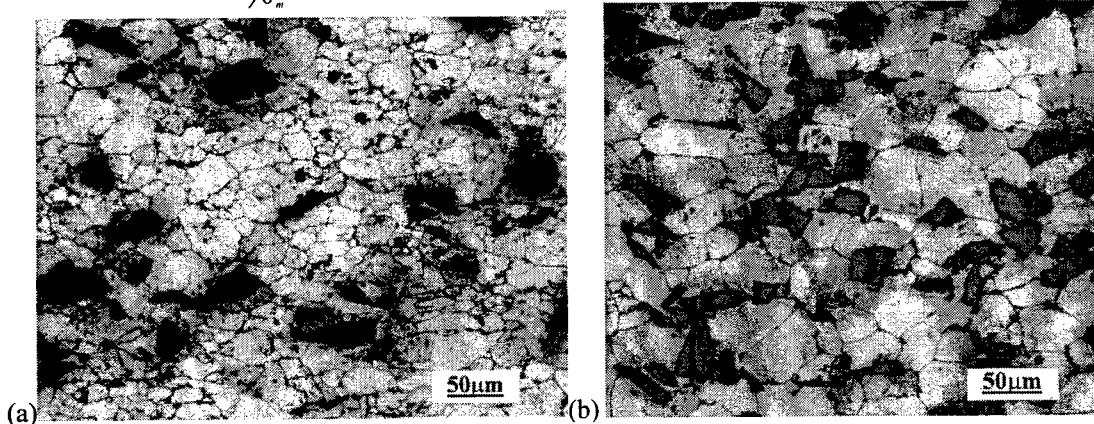


Fig. 7. SEM micrographs of HIPped MMCs containing (a) 5vol% and (b) 20 vol% SiCp showing particle cracking during tensile deformation.

CONCLUSION

Al-based MMCs were successfully fabricated by the HIPping process. The density of HIPped MMCs is close to the theoretical value, and is much higher than those fabricated by conventional PM sintering.

The Young's modulus of MMCs increases with increasing SiCp volume fraction. The Young's modulus calculation from the Shear-lag model is lower than the experimental data. On the other hand, the modified Eshelby model gives good prediction of the Young's modulus.

The fracture strength of MMCs decrease with increasing SiCp volume fraction and cracking of the SiCp particles upon tensile deformation was observed.

ACKNOWLEDGMENT

K. C. LAU acknowledges City University of Hong Kong for fellowship support for this research.

REFERENCES

1. M.P. Thomas, J.E. King, 1993. Effect Of Thermal And Mechanical Processing On Tensile Properties Of Powders Formed 2124 Al and 2124 Al-SiCp Metal Matrix Composites, *Mater. Sci. Technol.* 9: 742-753.
2. D.J. Lloyd, 1994. Particle Reinforced Aluminum And Magnesium Matrix Composites, *Int. Mater. Rev.* 39: 1.
3. A.J. Shakesheff, 1995. Aging and Toughness of Silicon Carbide Metal-Matrix Composites, *J. Mater. Sci.* 30: 2269-2276.
4. S. G. Song, N. Shi, G.T. Gray, J.A. Roberts, 1996. Reinforcement Shape Effects on The Fracture Behavior and Ductility of Particulate-Reinforced 6061-Al matrix Composites, *Metall. Mater. Trans. A*, 27: 3739-46.
5. Z.Y. Ma, Y.L. Li, Y. Liang, F. Zheng, J. Bi, S.C. Tjong, 1996. Nanometric Si₃N₄ Particulate-Reinforced Aluminum Composite, *Mater. Sci. Eng. A*, 219: 229-231.
6. Z.Y. Ma, S.C. Tjong, 1997. In-Situ Ceramic Particulate-Reinforced Al-Matrix Composites Fabricated by Reaction Pressing in TiO₂ (Titanium)-Al-B(B₂O₃) Systems, *Metall. Mater. Trans. A*, 28: 1931-42.
7. A. Niklas, L. Froyen, L. Delaey, L. Buekenhout, 1991. Comparative Evaluation Of Extrusion And Hot Isostatic Pressing As Fabrication Techniques For Al-SiC Composites, *Mater. Sci. Eng. A*, 135: 225-229.
8. E. Pagounis, M. Talvitie, V.K. Lindoos, 1996. Influence Of Reinforcement Volume Fraction And Size On The Microstructure And Abrasion Wear Resistance Of Hot Isostatic Pressed White Iron Matrix Composites, *Metall. Mater. Trans. A*, 27: 4171-4181.
9. J.M. Kunze, H.N. Wadley, 1997. The Densification Of Metal-Coated Fibers: Hot Isostatic Pressing Experiments, *Acta Mater.* 45: 1851-1865.
10. Y. Flom and R. J. Arsenault, 1989. *Acta Metall. Mater.*, 37: 2413-2423.
11. B. Wang, G.M. Janowski, B.R. Patterson, 1995. SiC Particulate Cracking In Powder Metallurgy Processed Aluminum Matrix Composite Materials, *Metall. Mater. Trans. A*, 26: 2457-2467.
12. B.Y. Lou, J.C. Huang, 1996. Failure Characteristics of 6061/Al₂O₃/15p and 2014/Al₂O₃/15p Composites as a Function of Loading Rate, *Metall. Mater. Trans. A*, 27: 3095-3107.
13. H.L. Cox, 1952. The Elasticity and Strength of Paper and Other Fibrous Materials, *J. Appl. Phys.* 3: 72-79.
14. M. Taya, R.J. Arsenault, 1989. *Metal Matrix Composites*, Pergamon Press, 25-28
15. M. Taya, R.J. Arsenault, 1987. A Comparison Between a Shear Lag Type Model and An Eshelby Type Model in Predicting the Mechanical Properties of a Short Fiber Composite, *Scripta Metall.* 21: 349-345.
16. V. C. Nardone, K. M. Prew, 1985. On the Strength of Discontinuous Silicon Carbide Reinforced Aluminum Composites, *Scripta Metall.* 20: 43-48.
17. T. Mura, 1987. Micro-mechanics of Defects in Solids, Martinis Nijhoff Publication.
18. M. Taya, R.J. Arsenault, 1989. *Metal Matrix Composites*, Pergamon Press 250-251
19. K.C. Lau, 1996. Characterization and mechanical Properties of SiC reinforced Aluminum Matrix Composites, Master of Philosophy Thesis, City University of Hong Kong.
20. T. Mori and K. Tanaka, 1973. Average Stress in Matrix and Average Elastic Energy of Materials with Misfitting Inclusions, *Acta Metall.* 21: 571-574.

Hydrostatic Extrusion of Composite Rod

Ui-Bin Tsai, Chi-Wei Wu, Ray-Quen Hsu

Department of Mechanical Engineering
National Chiao-Tung University
Hsin-Chu, Taiwan, R.O.C.

ABSTRACT

Hydrostatic extrusion is a process where the billet is completely immersed in pressurized liquid. Pressure at the die orifice is lower than that at the die entrance, and the resulting pressure differential causes metal flow toward the orifice or die exit. The force required to push the billet through the die is thus provided by the hydrostatic pressure instead of by a direct ram force.

Composite rods or wire are composed of two or more different materials, each material having its own distinctive mechanical characteristics. Because a composite rod billet deforms more uniformly in hydrostatic extrusion than in conventional extrusion, it is considered that hydrostatic extrusion is one of the most effective processes to manufacture composite rods and wires.

The aim of this work is to propose an analytic model for hydrostatic extrusion. In this model, the liquid surrounding the billet is considered as a kind of hydrodynamic lubrication. On the other hand, the upper-bound theorem is adopted to analyse the deformation of the composite rod or wire. A model describing plastic flow of these clad rods has been established. Work-hardening effects of the materials are taken into account in the model. The results show the deformation behavior of billets in hydrostatic extrusion under different processing conditions. Finally an experiment of hydrostatic extrusion of round rod has been conducted and its results found to be very close to the model analysis.

INTRODUCTION

Hydrostatic extrusion is one of the high pressure extrusion processes used in industry for manufacturing of composite rods. The isotropic pressure produced by hydrostatic extrusion increases ductility of the material, which in turn increases the forming ability of many "difficult-to-extrude" rod billet.

Avitzur[1], Yamaguchi[2] and Matsura[3] all conducted pioneering experiments on the hydrostatic extrusion of composite materials. Avitzur[4][5] and Osakada[6] proposed analytical models to predict different deformation patterns of a composite billet during hydrostatic extrusion. Wilson[7][8][9] investigated the lubrication characteristics of hydrostatic extrusion. Snidle[10][11] discussed the thermal effects of the lubrication film. Cho[12] combined hydrodynamic lubrication theory with the admissible velocity field proposed by Nagpal[13] for the hydrostatic extrusion of tubes.

Most of these studies dealt only with rod composed of one particular material, however, the constituent component of the composite clad rod billet is usually made of materials with distinctive characteristic, thus its extrusion is far more difficult to accomplish.

In this study, an upper bound method is integrated with hydrodynamic lubrication theory to analyse hydrostatic extrusion of composite rod. The energy dissipated both by deformation of the billet under extrusion and the lubrication film are considered, and the deformation pattern of the composite rod can thus be predicted by minimization of the total energy consumed.

HYDROSTATIC EXTRUSION MODEL

Fig. 1 shows an analytical model of the hydrostatic extrusion of composite rod. The thin shaded area indicates the lubrication film established under extrusion. The outer sleeve of the composite billet is not in

direct contact with the die surface because of the existence of this lubrication film. The final dimension of the product is determined by the thickness of this film, which in turn can be calculated from hydrodynamic lubrication theory proposed by Reynolds.

In hydrostatic extrusion, we made the following assumptions:

- (1) Both weight and inertia of the lubrication film are neglected.
- (2) The lubrication film is an incompressible Newtonian fluid.
- (3) The film thickness is small compared with the other dimensions.
- (4) No slip exists between fluid and the solid surfaces.

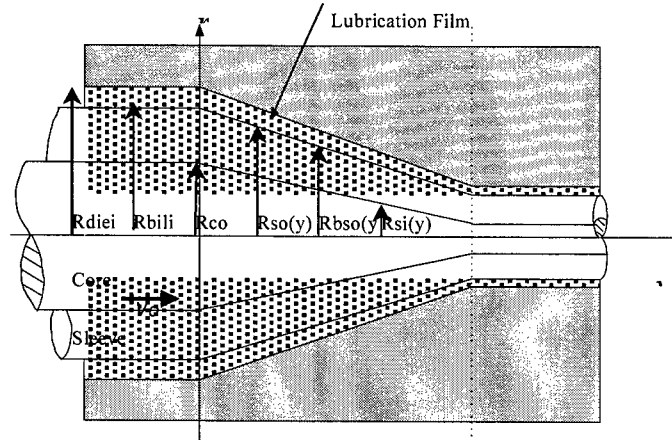


Fig. 1. Analytical model of the hydrostatic extrusion of composite rod.

Fig. 2. is a free body diagram of the lubrication film, from $\sum F_y = 0$, we have:

$$p \cdot r dr d\theta - (p + \frac{\partial p}{\partial y} dy) \cdot r dr d\theta - (\tau + \frac{\partial \tau}{\partial r} dr) \cdot (r + dr) d\theta dy + \tau \cdot r d\theta dy = 0 \quad 1.$$

here p is the pressure, τ is the shear force. Since $p = p(y)$? $\tau = \tau(r)$, Equation 1 can be reduced to:

$$\tau + r \cdot \frac{d\tau}{dr} = -\frac{dp}{dy} \cdot r \quad 2.$$

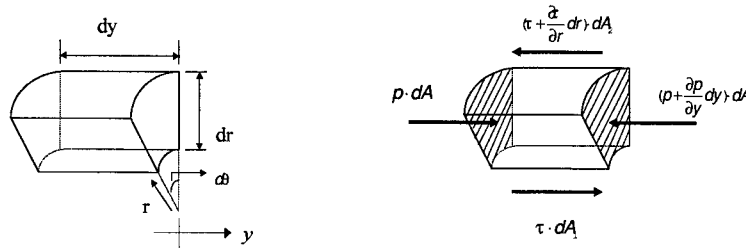


Fig. 2. Free body diagram of the lubrication film.

On the other hand, because the lubrication film is Newtonian fluid:

$$\tau = \mu \frac{du}{dr} \quad 3.$$

here, u is the velocity of the fluid. Replace (3) into (2), we now have:

$$\mu \cdot \frac{du}{dr} + r \cdot \frac{d}{dr} \left(\frac{du}{dr} \cdot \mu \right) = -\left(\frac{dp}{dy} \cdot r \right) \quad 4.$$

From assumption (4), we know there is no slip between fluid and the outer sleeve of the billet, so we determine on the outer sleeve surface that:

$$u = \frac{R_{bili}^2}{R_{bso}^2} \cdot V_0 \quad 5.$$

Also at the die surface $u = 0$, so:

$$u = \frac{1}{4\mu} \cdot \frac{dp}{dy} \left[(R_{so}^2 - r^2) - \frac{(R_{so}^2 - R_{bso}^2)}{\ln(\frac{R_{so}}{R_{bso}})} \cdot \ln(\frac{R_{so}}{r}) \right] + \frac{R_{bili}^2 \cdot V_0}{R_{bso}^2 \cdot \ln(\frac{R_{so}}{R_{bso}})} \cdot \ln(\frac{R_{so}}{r}) \quad 6.$$

By substituting Equation 6 into Equation 4, at constant flow rate, we can solve for dp/dy and with the help of boundary condition, we finally obtain the dissipation energy R consumed by the lubrication film per unit volume:

$$R = \mu \left(\frac{\partial u}{\partial r} \right)^2 \quad 7.$$

Total energy loss caused by the film \dot{W}_{lub} is:

$$\dot{W}_{lub} = 2\pi \int_0^L \int_{R_{bso}}^{R_{so}} R \cdot r dr dy \quad 8.$$

On the other hand, calculation of the extrusion energy is based on the upper bound method, the velocity fields of the composite billet is adopted from Jeng and Hsu[14], which under polar coordinate specify the following:

a. for the sleeve :

$$V_y(y) = \frac{V_0 \int_0^{\phi_f(0)} \{R_{so}^2(\phi, 0) - R_{si}^2(\phi, 0)\} d\phi}{\int_0^{\phi_f(y)} \{R_{so}^2(\phi, y) - R_{si}^2(\phi, y)\} d\phi} \quad 9.$$

$$V_\phi(r, \phi, y) = \frac{r}{R_{so}^2(\phi, y) - R_{si}^2(\phi, y)} \int_0^\phi \frac{\partial}{\partial y} \{ [R_{si}^2(\phi, y) - R_{so}^2(\phi, y)] \cdot V_y(y) \} d\phi \quad 10.$$

$$V_r(r, \phi, y) = - \left[\frac{r}{2} \left\{ \frac{\partial V_y(y)}{\partial y} + \frac{\partial \omega(\phi, y)}{\partial \phi} \right\} + \frac{R_{so}^2(\phi, y)}{2} \left\{ \frac{\partial V_y(y)}{\partial y} + \frac{\partial \omega(\phi, y)}{\partial \phi} \right\} \right. \\ \left. - \frac{1}{r} \left\{ R_{so}(\phi, y) \cdot \omega(\phi, y) \frac{\partial R_{so}(\phi, y)}{\partial \phi} + R_{so}(\phi, y) \cdot V_y(y) \frac{\partial R_{so}(\phi, y)}{\partial y} \right\} \right] \quad 11.$$

b. for the core :

$$V_y(y) = \frac{V_0 \int_0^{\phi_f(0)} R_{si}^2(\phi, 0) d\phi}{\int_0^{\phi_f(y)} R_{si}^2(\phi, y) d\phi} \quad 12.$$

$$V_\phi(r, \phi, y) = - \frac{r}{R_{si}^2(\phi, y)} \int_0^\phi \frac{\partial}{\partial y} \{ R_{si}^2(\phi, y) \cdot V_y(y) \} d\phi \quad 13.$$

$$V_r(r, \phi, y) = - \frac{r}{2} \left\{ \frac{\partial V_y(y)}{\partial y} + \frac{\partial \omega(\phi, y)}{\partial \phi} \right\} \quad 14.$$

Total energy \dot{J}_{total}^* consumed by the composite rod during hydrostatic extrusion thus equals,

$$\dot{J}_{total}^* = \dot{J}_{billet}^* + \dot{W}_{lub} \quad 15.$$

\dot{J}_{total}^* is directly coupled with the film thickness, so from the upper bound theory, we know that deformation of the billet will occur at the minimum energy consumed, in this way, the deformation pattern of the composite rod can be predicted.

RESULTS AND DISCUSSION

Common defects of the composite rod during extrusion includes fracture or wavy deformation of core or sleeve, non deformation of the core, etc. In this study, both core and sleeve are allowed to have their own exit velocity at the die exit. If the result indicates slip between core and sleeve at the die exit, we will conclude non-homogeneous deformation takes place, which will lead to product defect.

Fig. 3. shows the extrusion pressure with different semi-die angles for a Cu/Al(core) combination. For core volume percentages from 20~80, hydrostatic extrusion will render a sound deformation, while for the traditional extrusion only the unshaded area can be safely extruded[15]. This result clearly indicates hydrostatic extrusion is a much better forming process for manufacturing of the composite rod.

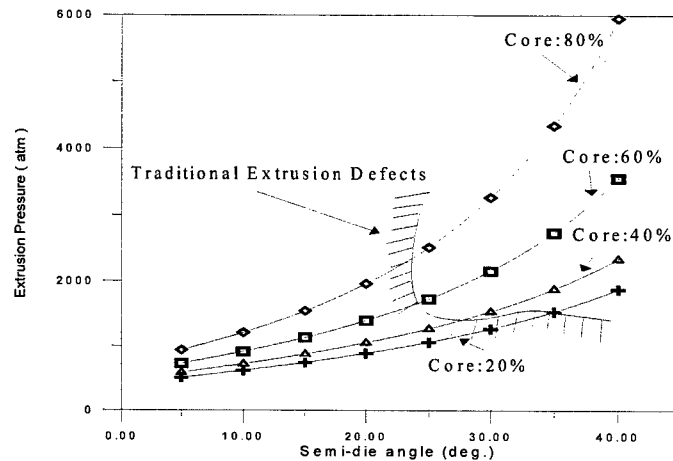


Fig. 3. Extrusion pressure with different semi-die angles for Cu/Al(core) combination.

Fig. 4. depicts extrusion pressure with the core/sleeve strength ratio. The semi-die angle is 10° and the core volume percentage is 64. For the various extrusion ratios shown, the composed rod can be safely extruded hydrostatically, while traditional extrusion can only be safely conducted at the unshaded area[15].

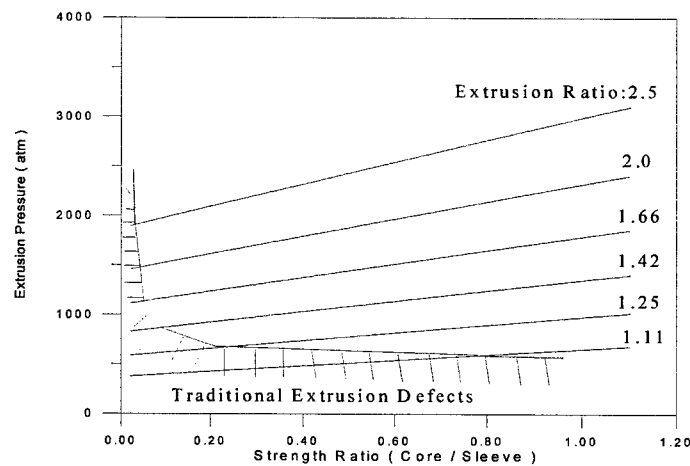


Fig. 4. Extrusion pressure with the strength ratio of core/sleeve.

Extrusion pressure increased with the increase of extrusion ratio is shown in Fig. 5. Here, the combination of the composite rod is O.F.H.C. Copper/Al-6061(core), semi-die angle is 15° . When compared with the experiment conducted by Yamaguchi[2], we found the results to be similar.

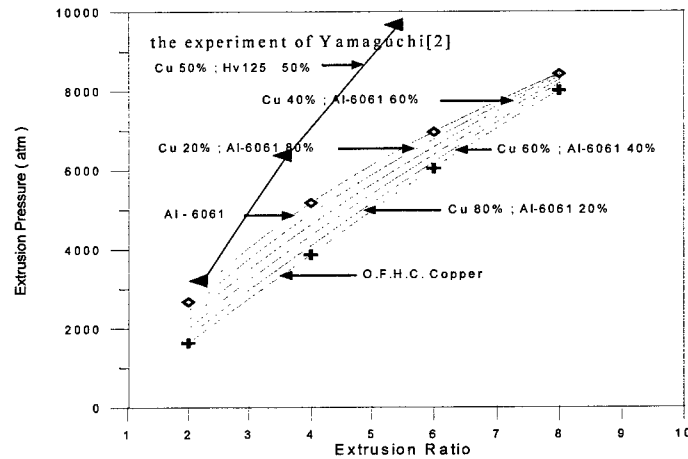


Fig. 5. Extrusion pressure vs. extrusion ratio.

The effect of strength ratio on the extrusion ratio is indicated in Fig. 6. When the extrusion ratio is 2.0 and the semi-die angle equals 15° , the larger the volume percentage of the stronger constituent, the larger the extrusion pressure.

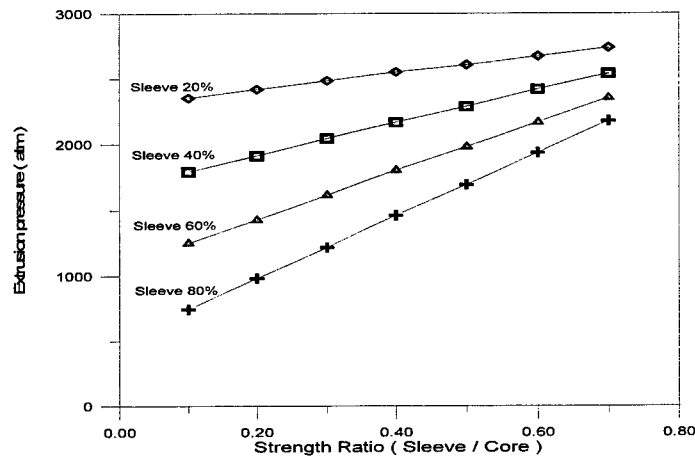


Fig. 6. Extrusion pressure versus strength ratio.

Finally, Fig. 7 shows the analytical results have the same trend as the experimental results. It is clear that this model proposed a numerical result higher than the experiment. This is because upper bound theory predicts a higher energy consumption than what is truly consumed.

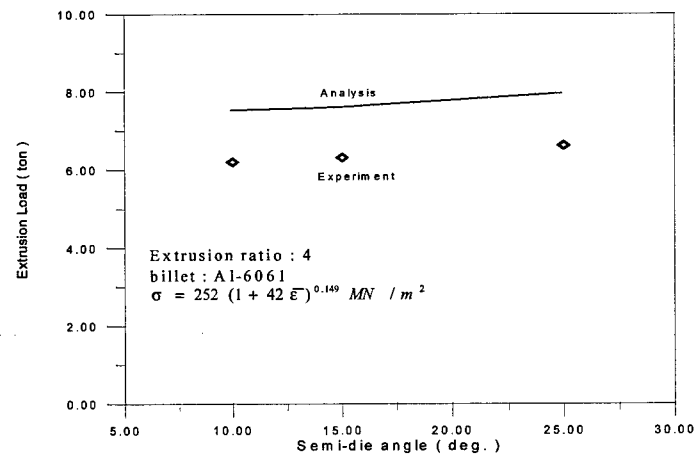


Fig. 7. Analytical results compared with experimental results.

CONCLUSION

An analytical model applicable to the hydrostatic extrusion of composite rod is proposed. It is found :

1. Hydrostatic extrusion has α far wider forming range than traditional extrusion.
2. Extrusion pressure increase with both extrusion ratio and strength ratio.
3. The combination of hydrodynamic lubrication theory with the upper bound theory obtained an analytical results agreeable with the experiment.

REFERENCE

1. Zoerner, W., Austen, A., Avitzur, B., 1972. Hydrostatic Extrusion of Hard Core Clad Rod. Trans. ASME , J. Eng. for Ind., 78-80.
2. Yamaguchi, Y., Noguchi, M., Matsushita, T., Nishihara, M., 1974. Hydrostatic Extrusion of Clad Materials. J. Japan Society for Technology of Plasticity, 15(164) , 723-729.
3. Matsuura, Y., Takase, K., 1974. An Experimental Study and Solution of the Energy Method on Plastic Deformation of Two-Phase Combination Materials Consisting of Copper and Aluminum - study of characteristics on hydrostatic extrusion of combination materials I., ? ? ? ? , 15 (157) , 156-165.
4. Avitzur B., 1965. Hydrostatic Extrusion. Trans. ASME, 87, 487-494.
5. Avitzur, B., 1972. Experiment Study of Hydrostatic Extrusion. Trans. ASME , 658-668.
6. Osakada, K., Limb, M., Mellor, P.B., 1973. Hydrostatic Extrusion of Composite Rods with Hard Cores. Int. J. Mech. Sci., 15, 291-307.
7. Wilson, W.R.D., 1974. A Thermal Reynolds Equation and Its Application in the Analysis of Plasto-Hydrodynamic Inlet Zones. J. Lubrication Technology, Trans. ASME , 95(4), 572-577.
8. Wilson, W.R.D., 1973. The Variation of Lubrication Film Thickness in the Work Zone of Hydrodynamically Lubricated Continuous Deformation Processes. J. Lubrication Technology , Trans. ASME , 95(4), 541-543.
9. Wilson, W.R.D. , 1976. Hydrodynamic Lubrication of Hydrostatic Extrusion. J. Lubrication Technology , Trans. ASME, 541-543.
10. Snidle R.W., 1973. An Elasto-Plasto-Hydrodynamic Lubrication Analysis of the Hydrostatic Extrusion Process., J. Lubrication Technology, Trans. ASME , 95(2) , 113-122.
11. Snidle, R.W. A., 1976. Thermal Hydrodynamic Lubrication Theory for Hydrostatic Extrusion of Low Strength Materials. J. Lubrication Technology, Trans. ASME, 98 , 335-343.
12. Cho, N.S., 1983. Hydrofilm Extrusion of Tubes Through Optimized Curved Dies. Trans. ASME , J. Eng. for Ind , 105, 243-250.
13. Nagpal, V., 1974. General Kinematically Admissible Velocity Fields for Some Axisymmetric Metal Forming Problems. Trans. ASME, J. Eng. for Ind., 1197-1201.
14. Jeng, J.L., Hsu, R.-Q., 1996. Extrusion of Composite Multi-Core Clad Rods Composed of Three and More Different Materials. Master Thesis , National Chiao-Tung University.
15. Avitzur, B., 1982. Criterion for the Prevention of Core Fracture During Extrusion of Bimetal Rods. Trans. ASME , J. Eng. for Ind , 104, 93-304.

Numerical Modelling and Localized Failure Analysis in Metal Powder Forming Processes

A.R. Khoei, R.W. Lewis, D.T. Gethin

Mechanical Engineering Department, University of Wales Swansea,
Singleton Park, Swansea, SA2 8PP, U.K

ABSTRACT

In this paper, a general framework for the finite element simulation of powder forming processes is presented. The research is focused on two different operations: the compaction process of powder and the process of localization. In the process of compaction, powders exhibit strain or work hardening, the volume reduces and the material becomes harder. A model is adopted to represent this physical process by employing an updated Lagrangian formulation for large deformation, a hardening cap model for the compressible behaviour of powder material and an interface element formulation for frictional behaviour of the contact surface. An adaptive analysis based on error estimates and automatic remeshing techniques is also applied to simulate the compaction process. In the process of localization, the ultimate capacity of the new materials is evaluated using the analysis of failure. A method for dealing with incompressible material is presented for capturing the strain localization and displacement discontinuity. An adaptive procedure is also introduced for elongating elements with suitable error measures, which will indicate and capture effectively the localization regions.

INTRODUCTION

The numerical simulation of the compaction process is central to an understanding of the mechanics of powder behaviour and when it is coupled with experimental inputs the simulation can be considered as an alternative tool to achieve a more economic enterprise. The computational models developed for the modelling of compaction forming processes can be basically classified into two approaches, i.e. 'micro-mechanical' and 'macro-mechanical' approaches. The micro-mechanical approach considers the discrete nature of powder particles, whereas the macro-mechanical, or continuum approach, characterises the overall behaviour of the powder mass by idealising the powder mass as an equivalent continuum material. Both the *micro-* and *macro-mechanical* approaches, have advantages and disadvantages in the modelling of powder compaction. Nonetheless, from an industrial viewpoint the macro-mechanical approach has a definitive edge over the micro-mechanical approach in that the gross behaviour of the powder mass can be modelled and simulated on an industrial scale. In the present study a finite element model based on the macro-mechanical approach is adopted to present two different operations: the compaction process of powder and the process of localization.

The compaction forming of metal powder is a process involving large deformations, large strain, non-linear material behaviour and friction. For a successful modelling of such a highly non-linear behaviour, a cap plasticity model is applied to describe the constitutive model of compressible and hardening behaviour of materials [1]. This model reflects the yielding, frictional and densification characteristics of powder along with strain and geometrical hardening which occurs during the compaction process. A hardening rule is used to define the dependence of the yield surface on the degree of plastic straining. A plasticity theory for friction is employed in the treatment of the powder-tooling interface [2]. The involvement of two different materials, which have contact and relative movement in relation to each other, must be considered. A special formulation for friction modelling is coupled with a material formulation. The interface behaviour between the die and powder is modelled by using an interface element mesh.

The process of localization refers to the phenomenon by which the deformation in solids localize into narrow bands of intense straining. It is well known that strain localization and indeed displacement discontinuity can arise in materials exhibiting plastic behaviour. Indeed such localization is almost certain

to occur if strain softening or non-associated behaviour exists, though it can be triggered even when ideal plasticity is assumed. This study is concerned mainly with the manner in which the numerical discretization process has to be devised so as to capture the localization phenomenon.

THE MIXED $\mathbf{u} - \pi$ FINITE ELEMENT FORMULATION

The main problem in the application of a standard displacement formulation to incompressible, or nearly incompressible problems, lies in the determination of the mean stress, or pressure, which is related to the volumetric part of the strain. For this reason it is convenient to separate this from the total stress field and treat it as an independent variable. In the present study, a method is presented for applying the mixed formulation for both compressible and incompressible material [3]. The mixed $\mathbf{u} - \pi$ formulation for elasto-plastic analysis can be presented as

$$\begin{aligned} \mathbf{S}^T(\boldsymbol{\sigma}_d - \mathbf{m}\pi) + \rho \mathbf{b} - \rho \ddot{\mathbf{u}} &= 0 \\ \mathbf{m}^T \boldsymbol{\varepsilon} - \frac{\pi}{K} &= 0 \end{aligned} \quad 1.$$

where \mathbf{u} is the displacement vector, \mathbf{b} is the body force acceleration, ρ is the density and $\boldsymbol{\sigma}$ is the total stress. $\boldsymbol{\sigma}_d$ is the deviatoric stress vector defined as $\boldsymbol{\sigma}_d = \boldsymbol{\sigma} - \mathbf{m}\pi$, the mean stress π is $\pi = (1/3) \mathbf{m}^T \boldsymbol{\sigma}$ with \mathbf{m} denoting a vector which has a form of $\mathbf{m}^T = [1, 1, 1, 0, 0, 0]$ for the general three-dimensional stresses. \mathbf{S} is the strain operator relating displacements and strain ($\boldsymbol{\varepsilon} = \mathbf{S}\mathbf{u}$), K is the bulk modulus of material defined as $K = E / [3(1-2\nu)]$, where E is the elastic Young's modulus and ν is the Poisson's ratio. Applying the standard finite element Galerkin discretization process to equations (1), with the independent approximations of \mathbf{u} and π defined as, $\mathbf{u} = \mathbf{N}_u \bar{\mathbf{u}}$ and $\pi = \mathbf{N}_\pi \bar{\pi}$, we obtain the following algebraic equations

$$\begin{aligned} \int_{\Omega} \mathbf{B}^T \boldsymbol{\sigma}_d d\Omega + \mathbf{Q} \bar{\pi} + \mathbf{M} \ddot{\bar{\mathbf{u}}} &= \mathbf{f}_u \\ \mathbf{Q}^T \bar{\mathbf{u}} - \mathbf{C} \bar{\pi} &= \mathbf{f}_\pi \end{aligned} \quad 2.$$

where \mathbf{B} is the strain matrix relating the increments of strain and displacement (i.e. $d\boldsymbol{\varepsilon} = \mathbf{B} d\bar{\mathbf{u}}$). If displacements are large, then the strains depend in a non-linear manner on the displacements [2]. \mathbf{M} , \mathbf{Q} and \mathbf{C} are the mass, coupling and compressibility matrices, defined as

$$\mathbf{M} = \int_{\Omega} \mathbf{N}_u^T \rho \mathbf{N}_u d\Omega, \quad \mathbf{C} = \int_{\Omega} \mathbf{N}_\pi^T \frac{1}{K} \mathbf{N}_\pi d\Omega \quad \text{and} \quad \mathbf{Q} = \int_{\Omega} \mathbf{B}^T \mathbf{m} \mathbf{N}_\pi d\Omega \quad 3.$$

and \mathbf{f}_u and \mathbf{f}_π vectors are defined as

$$\mathbf{f}_u = \int_{\Omega} \mathbf{N}_u^T \rho \mathbf{b} d\Omega + \int_{\Gamma} \mathbf{N}_u^T \mathbf{t} d\Gamma \quad \text{and} \quad \mathbf{f}_\pi = 0 \quad 4.$$

The definition of the spatial problem is complete when the constitutive law for use in the first term of equation (2) is defined. This term represents the internal force, and for non-linear problems, can be written as

$$\int_{\Omega} \mathbf{B}^T \boldsymbol{\sigma}_d d\Omega = \int_{\Omega} \mathbf{B}^T (\mathbf{D}_T - \mathbf{m} K \mathbf{m}^T) \mathbf{B} \bar{\mathbf{u}} d\Omega = \mathbf{K}_T \bar{\mathbf{u}} \quad 5.$$

where \mathbf{D}_T is the consistent tangential stiffness matrix can be obtained by performing a full differentiation on the internal force term.

THE COMPACTION OF MATERIAL

In the mechanical behaviour of metal powder forming, the focus of attention has been mainly on two points, the choice of suitable constitutive models and the frictional algorithm for describing the interaction at the powder-die interface. For the nonlinear behaviour of powder materials, a double-surface plasticity model based on a combination of a convex yield surface consisting of a failure envelope, such as a Mohr-Coulomb yield surface and a hardening elliptical yield cap is developed [1]. This model reflects the

yielding, frictional and densification characteristics of powder along with strain and geometrical hardening which occur during the compaction process. The solution yields details on the powder displacement from which it is possible to establish the stress state in the powder and the densification can be derived from consideration of the elemental volumetric strain.

For the interface behaviour between powder and die wall, a plasticity theory of friction which is similar to the theory of elasto-plasticity based on a stick (or adhesion) law, a stick-slip law, a wear and tear rule, a slip criterion and a slip rule, is employed [2]. The constitutive modelling of the frictional behaviour of the metal powder is modelled by Coulomb's friction law and the plasticity theory of friction in the context of an interface element formulation. Further details about the considered cap plasticity model for powder materials and friction model for interface between the die and powder are given in Refs.[1, 2].

ERROR ESTIMATION AND ADAPTIVE REMESHING

Despite the advantage of the Lagrangian description in which deformation history dependent variables (strain hardening) can easily be taken into account and the changing shape of the formed product can be followed, one of the significant limitations of this approach is the progressive distortion. This can cause large approximation errors or make the Jacobian determinant negative. In order to solve this problem, an efficient way is to use error estimates and adaptive remeshing, which not only control the discretization error, but automates the simulation process [4, 5]. The error estimates, which we shall use here, will be presented in the L_2 norm of strain as

$$\|\mathbf{e}_\epsilon\| = \|\epsilon^* - \hat{\epsilon}\| \quad 6.$$

where $\hat{\epsilon}$ is the discontinuous strain derived by a finite element solution and ϵ^* represent an improved solution which can be obtained by a global smoothing using the interpolation function N_u and nodal parameter $\bar{\epsilon}^*$ as

$$\epsilon^* = N_u \bar{\epsilon}^* \quad \text{where} \quad \bar{\epsilon}^* = \left(\int_{\Omega} N_u^T N_u d\Omega \right)^{-1} \int_{\Omega} N_u^T \hat{\epsilon} d\Omega \quad 7.$$

The adaptive mesh refinement strategy depends on the nature of the criteria on accuracy which we wish to satisfy. A very common requirement is to specify the achievement of a certain minimum percentage error in the L_2 norm. Thus, we require that after remeshing each element will have the same error and the overall percentage error is equal to some target percentage error (i.e. $\eta \leq \eta_{aim}$). If we assume that the error is equally distributed between elements these requirements can be translated into our placing a limit on the error in each element and the new element size can be predicted.

A Rotational Flanged Component

In order to illustrate the applicability of the present formulation, an adaptive FEM analysis has been employed for the non-linear behaviour of a rotational flanged component, as illustrated in Figure 1(a). The analysis commences with an uniform distribution of elements. The non-linear solution is carried out for various loading steps until the estimated error exceeds a prescribed value, or the sign of the Jacobian determinant becomes negative. At this time a new mesh is generated using appropriate refinement criterion. The information at the end of the previous step is transferred to the integration points of the new mesh by using nodal points as reference points to store the information. Finally, the non-linear procedure is performed on the new mesh starting from the beginning of the new step.

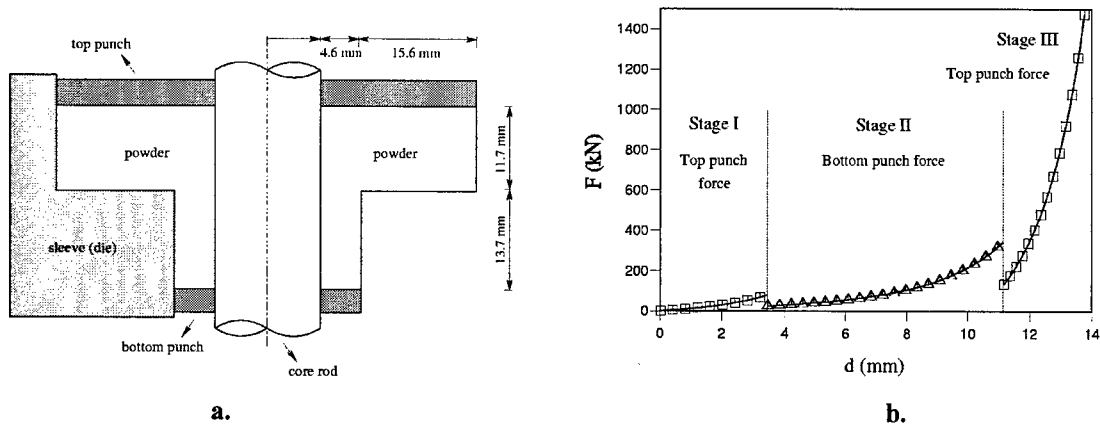


Fig. 1. A rotational flanged component; a) Geometry and boundary conditions, b) Predicted top and bottom punches forces at different displacements.

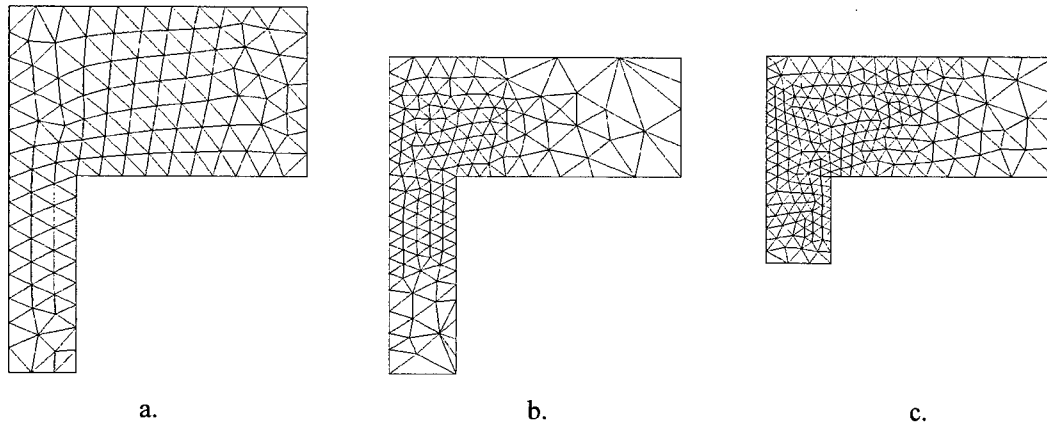


Fig. 2. Adaptive FEM analysis for a rotational flanged component; a) Initial mesh, b) First adaptive mesh ($\eta_{aim} = \% 4$), c) Second adaptive mesh ($\eta_{aim} = \% 2$).

The iron powder properties chosen in the analysis are given in Ref.[5]. The variations of top and bottom punch forces with displacements are plotted in Figure 1(b). The compaction process is carried out by a top punch movement of 3.44 mm, then a bottom punch movement of 7.70 mm and finally a further top punch movement of 2.62 mm, as illustrated in Figure 2. It can be seen that the proposed adaptive finite element approach is capable of simulating metal powder compaction processes in an efficient and accurate manner.

THE PROCESS OF LOCALIZATION

One of the other major challenges in the computation of powder forming problems is the analysis of failure. Such computations are needed to evaluate the ultimate capacity of new materials, but they are fraught with serious difficulties. One of these difficulties is the process of localization, which is ubiquitous in failure. In this study, the mixed $u - \pi$ formulation is applied for incompressible plasticity material such as a Tresca or Von-Mises material. In particular, if plasticity of an isochoric (volume-preserving) type is used, typified by classical Tresca or Von-Mises models, difficulties will arise with a simple, linear triangle (3C) finite element [3]. These difficulties are avoided to some extent by using the equivalence of the displacement form with a mixed formulation involving both the displacement u and pressure π as independent variables. It was shown in Ref.[3] that if a correct approximation is used then both the uniform and non-uniform mesh refinements will converge to the correct answer and clearly indicate the localization phenomenon.

In order to evaluate the ultimate capacity of the new compacted powder materials, a flanged component, at the final stage of compaction (Figure 2c), is numerically analysed by applying the mixed $u - \pi$ formulation

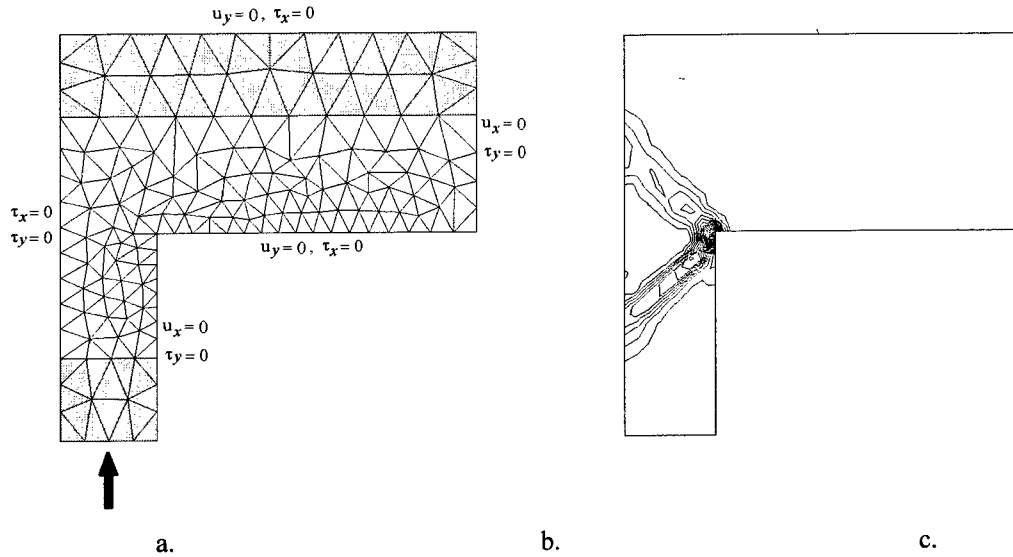


Fig. 3. A compacted flanged component; a) Finite element modelling using T6C/3C elements, b) Effective plastic strain contour at $d = 1.12$ mm ($H = -5000$ N/m²).

and using the triangular quadratic continuous displacement elements with the triangular continuous linear pressure elements (T6C/3C). A single movement of the bottom punch is applied to capture the strain softening. The finite element modelling of this compacted flanged component along with the effective plastic strain contour at $d = 1.12$ mm are given in Figure 3. The variation with displacement of the reaction of the bottom punch for two values of the plastic hardening/softening modulus, i.e. $H=0.0$ for ideal plasticity and $H=-5000$ N/m² for softening plasticity are plotted in Figure 4(a). It can be seen that the mixed formulation can be effectively used for such a compressible-incompressible combined material. However, with finite elements the problem of the local element size influencing the final solution for strain softening materials remains and an adaptive analysis using element elongation can be effective in the modelling of such phenomena.

Adaptive Analysis Using Element Elongation

An adaptive analysis is presented for elongating elements with suitable error measures, which will indicate and capture effectively the localization regions. The h -adaptive remeshing procedure adopted here is based on the manner developed by Zienkiewicz *et al.*[6]. In order to measure and indicate the error occurring in the values of displacements in each individual element, we shall estimate its magnitude based on the gradient of displacements as

$$e = \|u\| - \|u^h\| \leq \tilde{e} = ch \left| \frac{\partial u^h}{\partial \bar{x}} \right|_{\max} \quad 8.$$

where $\|u\|$ and $\|u^h\|$ are the exact and finite element solutions, \bar{x} is any direction chosen, h is elemental size and c is a positive constant. In practical application we shall try to devise a procedure in which the error indicator in each element is reduced to, or below, a specified values. In this case, a solution including displacements and their derivatives will be obtained from an initial mesh. The minimum size of elements and their directions will be indicated using expression (8). Finally, a mesh satisfying the requirements will be generated by the best available procedure.

An adaptive analysis using element elongation is performed to the compacted flanged component, given in Figure 3, at $d = 1.12$ mm from the initial mesh (Figure 3a) in the manner described above. As expected the adaptive remeshing provides a clear picture of localization phenomena, as illustrated in Figure 4(b). The numerical analysis shows that adaptive solutions using elongating element can be effective whatever starting mesh is used.

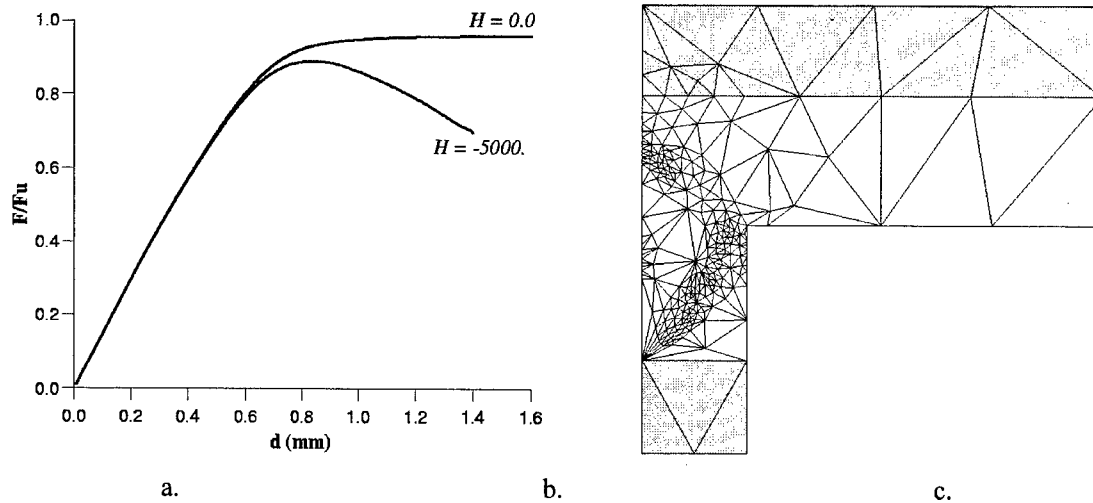


Fig. 4. a) Reaction of the bottom punch, b) Adaptive analysis using element elongation.

CONCLUSION

In this paper, a numerical analysis of metal powder forming processes was presented for both compressible and incompressible materials. In the process of compaction, an adaptive analysis was simulated by an updated Lagrangian finite element formulation. For the adaptive strategy, a posteriori Zienkiewicz-Zhu estimator using L_2 norm of strain by a recovery procedure was proposed. To describe the constitutive model of compressible behaviour of powder materials, a cap plasticity model using a hardening rule was applied to define the dependence of the yield surface on the degree of plastic straining. A special formulation for friction modelling was coupled with a material formulation by employing the interface element mesh in the contact area between the die and powder. In order to demonstrate a part of the wide range of problems that can be solved by the present formulation, the powder behaviour during the compaction of a rotational flanged component was analysed numerically. In a further work, the ultimate capacity of new materials was evaluated using the analysis of failure. A mixed $u - \pi$ formulation was applied for incompressible plasticity material such as a Tresca or Von-Mises material. An adaptive analysis using element elongation was performed to the compacted flanged component. The numerical analysis shows that adaptive solutions using elongating element can be effective whatever starting mesh is used. The results clearly indicate that the algorithm makes it possible to simulate the powder forming problems efficiently and automatically.

REFERENCES

1. A.R. Khoei, R.W. Lewis, 1998. 'Finite element simulation for dynamic large elasto-plastic deformation in metal powder forming', *Finite Elements in Analysis and Design*, 30, 335-352.
2. R.W. Lewis, A.R. Khoei, 1998. 'Numerical modelling of large deformation in metal powder forming', *Computer Methods in Applied Mechanics and Engineering*, 159, 291-328.
3. A.R. Khoei, R.W. Lewis, O.C. Zienkiewicz, 1997. 'Application of the finite element method for localized failure analysis in dynamic loading', *Finite Elements in Analysis and Design*, 27, 121-131.
4. O.C. Zienkiewicz, G.C. Huang, Y.C. Liu, 1990. 'Adaptive FEM computation of forming processes, Application to porous and non-porous materials', *Inter. J. Num. Meth. Eng.*, 30, 1527-1553.
5. A.R. Khoei, R.W. Lewis, 1999. 'Adaptive finite element remeshing in a large deformation analysis of metal powder forming', *Inter. J. for Numerical Methods in Engineering* (in press).
6. O.C. Zienkiewicz, M. Pastor, M. Huang, 1995. 'Softening, localization and adaptive remeshing. Capture of discontinuous solutions', *Computational Mechanics*, 17, 98-106.

Microstructure and High Temperature Deformation Behavior of a Tin / Ti_5Si_3 Nano-Grain Composite Produced by Non-Equilibrium PM Process

Kei Ameyama* and Yasuhiko Suehiro**

*Department of Mechanical Engineering,
College of Science and Engineering,
Ritsumeikan University,
1-1-1 Noji-Higashi, Kusatsu city, Shiga 525-8577, Japan
** Graduate Student, Ritsumeikan University

ABSTRACT

Microstructure and high temperature deformation behavior of Ti- Si_3N_4 mechanically alloyed (MA) powder compacts were investigated. Powders of the elements Ti and Si_3N_4 whose composition was Ti-20mass% Si_3N_4 were blended for MA. A planetary ball mill was used for milling under an Ar gas atmosphere. The MA powder milled for 720ks was heated at various temperatures. The MA powder was consolidated by vacuum hot press (VHP) at 200MPa for 10.8ks at 803K. The specimens were provided for compression tests at 913K~1073K at various initial strain rates. The MA process of Ti and Si_3N_4 powders for 720ks resulted in the formation of an amorphous and an α -Ti phases. These phases changed to TiN, Ti_2N and Ti_5Si_3 phases after the heat treatment at elevated temperatures. A (TiN + Ti_5Si_3) ultra fine microduplex structure was obtained after the heat treatment at 1473K for 3.6 ks. The compression tests revealed that the 803K-VHP specimen with non-equilibrium phases show the lowest flow stress at 993K at initial strain rate of $4.2 \times 10^{-4}\text{s}^{-1}$ in the three VHP specimens. Furthermore, the 803K- VHP specimen indicated lower flow stress at 993K at an initial strain rate of $4.2 \times 10^{-4}\text{s}^{-1}$ rather than that at $2.1 \times 10^{-4}\text{s}^{-1}$. Such a reverse of the flow stress between two different initial strain rates was attributed to the phase transformation during the deformation. The slower strain rate test produced larger amount of harder phases such as TiN, Ti_2N and Ti_5Si_3 . The specimen compressed to 25% ($\epsilon = 0.28$) at 993K at an initial strain rate of $4.2 \times 10^{-4}\text{s}^{-1}$ consisted of an (α -Ti+ Ti_2N + Ti_5Si_3) microduplex structure with an average grain size of approximately 40 nm. Therefore, there exists an appropriate condition for a low temperature and high strain rate forming process. A (TiN + Ti_5Si_3) microduplex structure with an average grain size of approximately 250 nm was also obtained in the specimen compressed to 25% after annealing at 1473K for 3.6 ks.

INTRODUCTION

The obstructions such as poor ductility and toughness are serious problem to be overcome for structural ceramics. The composite techniques originating from cermets or carbides have been developed to various types of composites and their processing technique [1]. Mechanical alloying (MA) was originally invented as a means of producing fine and uniform dispersion of oxides in superalloy matrix, and has a high potential of fabricating new types of composites. MA is one of the advantageous powder metallurgy (PM) processing technique which enables to give a non-equilibrium state to powders by introducing high density strain energy. In other words, the non-equilibrium PM processing such as MA is a kind of thermomechanical treatment of powders which can control microstructure of the sintered compacts. Application of the MA process to the titanium aluminides made a remarkable improvement on their ductility and workability by grain size refinement [2, 3]. It is, therefore, very important to pay an attention to the conditions of the treatment, such as heating rates and holding temperatures during sintering, as well as microstructure of the powders to obtain compacts with an ultra-fine grain structure.

EXPERIMENTAL PROCEDURE

The starting materials were a commercially pure Ti (99.7 mass%) powder of average particle size of 45 μm and an α - Si_3N_4 (98.7 mass%) powder of average particle size of 720 nm. They were blended to the composition of a Ti-20 mass% Si_3N_4 , and mechanically alloyed by a planetary ball mill conducted at a rotation speed of

250 rpm for 720 ks with a SKD11 vial and SUS304 stainless steel balls under an Ar-gas atmosphere. The powder to ball weight ratio was 1 : 3.6. During milling, 0.5 mass% of n-Heptane was added as a process control agent.

The 720ks MA powder was vacuum hot pressed (VHP) at 803 K for 10.8 ks, under a pressure of 200 MPa and at a heating rate of 0.5 K s^{-1} , followed by furnace cooling. The compression-test specimens of the column of $3 \text{ mm } \phi \times 4 \text{ mm}$ were cut from the VHP compacts. The VHP compacts were compression tested at temperature range between 913 K and 1073 K at various initial strain rates between $2.1 \times 10^{-4} \text{ s}^{-1}$ and $2.1 \times 10^{-3} \text{ s}^{-1}$. The specimens were examined by DSC, X-ray diffraction (XRD), SEM and TEM.

RESULTS AND DISCUSSION

The microstructure change during milling was examined by means of XRD analysis of the MA powders. Figure 1 shows XRD patterns using Cu Ka radiation of the powder mechanically alloyed for 0 s, 72 ks, 180 ks, 360 ks and 720 ks. Substantial broadening of the XRD peaks of the original species took place with the progress of milling, and XRD peaks of β -Ti phase was the only remaining peaks after 180 ks milling.

Figure 2 shows a TEM micrograph and a selected area diffraction pattern (SADP) of the MA powder milled for 720 ks. Although the SADP indicates that the MA powder is consisted mainly of an amorphous powder, a very fine grain was observed inside of the MA powder. TEM/EDS analysis revealed that such a fine grain had Ti richer composition rather than the matrix phase, so that the Ti rich nano-grain might be an α -Ti which is remained after MA process. Therefore, the MA powder provided for VHP compaction had a microstructure of an amorphous phase with a small amount of Ti-rich nano-grains.

The DSC analysis performed on the MA powder seem to support the actual existence of an amorphous phase. The DSC result of the 720 ks MA powder, shown in Figure 3, indicates that a large exothermic peak at 953 K took place due to the crystallization of the remaining amorphous phase in the MA powder. XRD examination of the MA powder heat treated at 973 K confirmed the formation of α -Ti, TiN, Ti_2N and Ti_5Si_3 phases. In addition, no carbide caused by the presence of n-Heptane was observed.

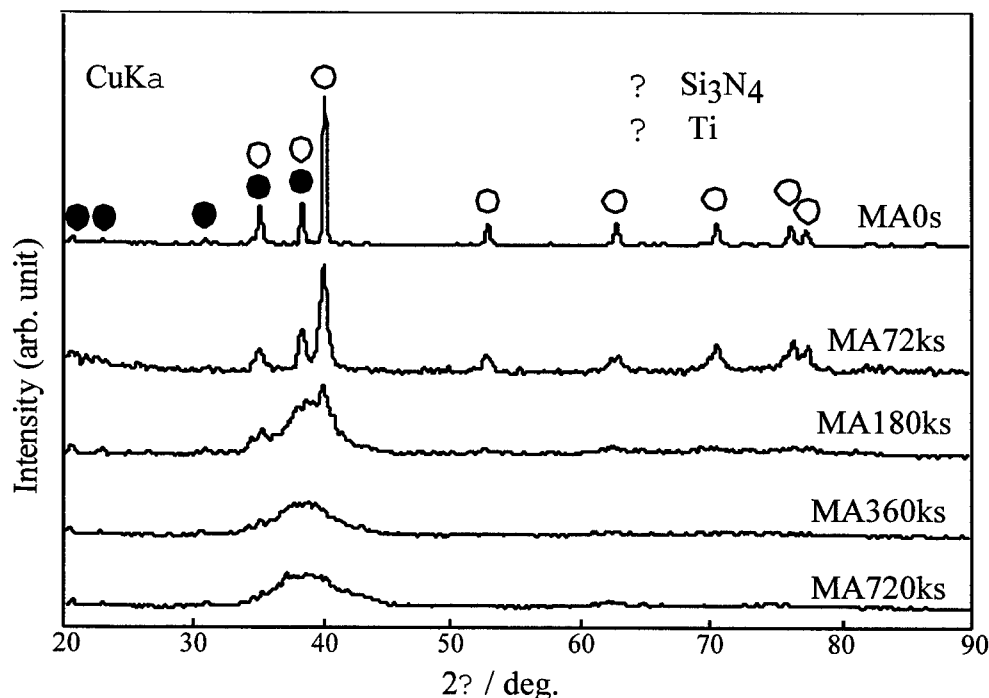


Fig. 1. XRD patterns of the MA powder milled for 0 s, 72 ks, 180 ks, 360 ks and 720 ks.

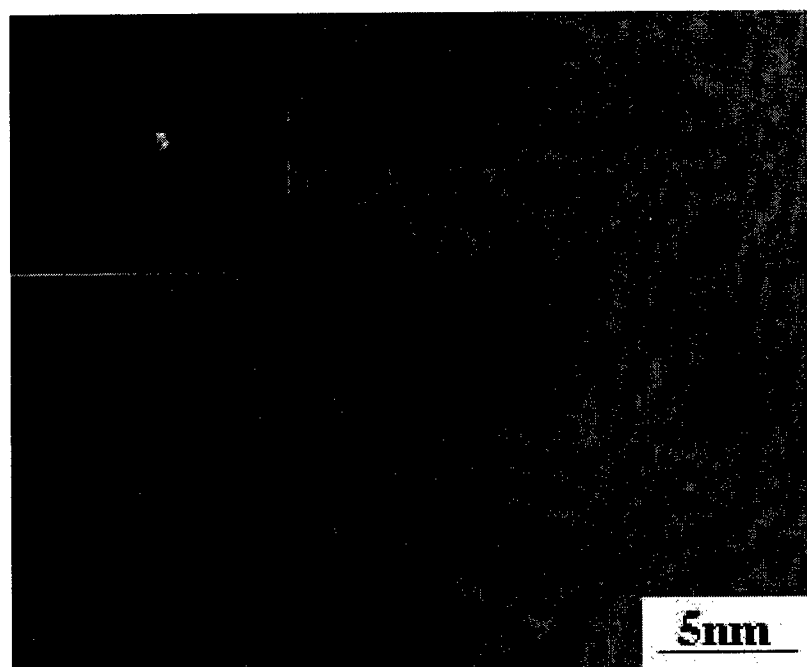


Fig. 2. TEM micrograph and a selected area diffraction pattern (SADP) of the MA powder milled for 720 ks.

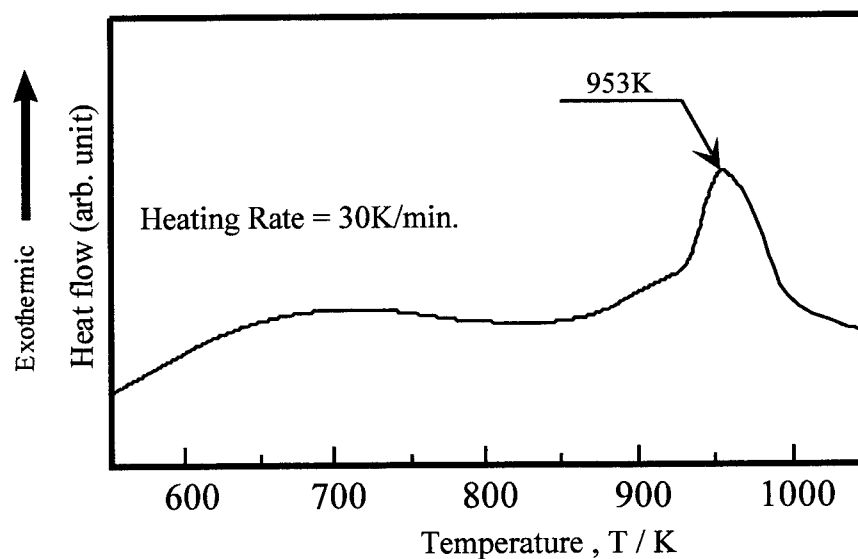


Fig. 3. DSC profile of the MA powder milled for 720 ks.

From these results, vacuum hot pressing was carried out at 803 K, in which the MA compact was expected to have an α -Ti nano-grain structure. Figures 4 (a) and (b) shows a XRD pattern of the VHP compact and a stress-strain curve for the compression test at 993 K at an initial strain rate of $2.1 \times 10^{-4} \text{ s}^{-1}$. Prior to each compression test the specimen was held at 993 K for 0.6 ks. The compression test was interrupted at a strain of 0.28. As can be seen in Figs. 4 (a) and (b), although the VHP compact was consist of an α -Ti phase at first, flow stress increased with strain during the compression test. In general, flow stress at elevated temperature decreases with test temperature. Therefore, such an unusual high temperature deformation behavior implies that structural change occurred with the compression test.

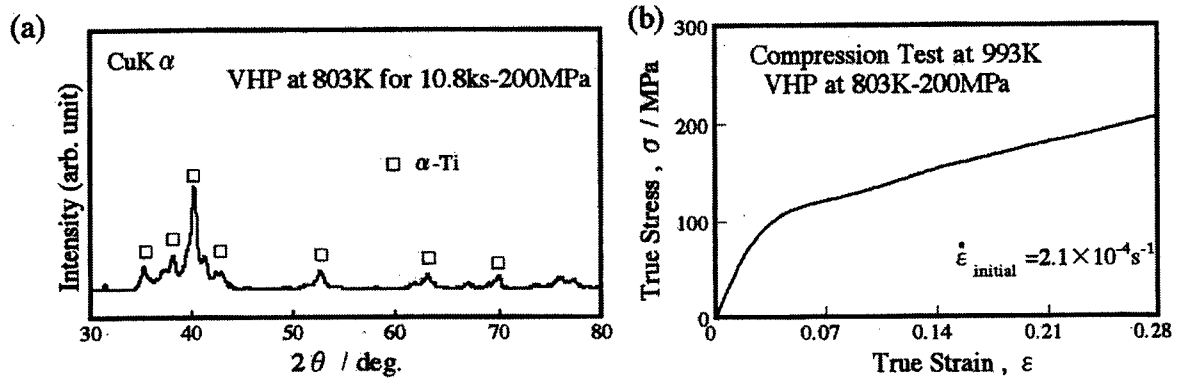


Fig. 4. a. XRD pattern of the VHP compact and
b. stress-strain curve for the compression test
at 993 K at an initial strain rate of $2.1 \times 10^{-4} \text{ s}^{-1}$.

Figure 5 shows XRD results of the VHP compact deformed to $\epsilon =$ (a) 0.00, (b) 0.06 and (c) 0.28 at 993 K. The VHP compact contains a large amount of α -Ti phase and a small amount of Ti_2N and Ti_5Si_3 phases at the initial state, i.e., $\epsilon = 0.00$. The latter two phases seem to be formed by nucleation and growth during holding at 993 K for 0.6 ks. However, as deformation at 993 K proceeds, the XRD peak intensity of α -Ti phase decreases while that of Ti_2N and Ti_5Si_3 phases increases. Figures 6 (a), (b) and (c) are the TEM micrographs of the specimen deformed to $\epsilon =$ (a) 0.00, (b) 0.06 and (c) 0.28 at 993 K at an initial strain rate of $2.1 \times 10^{-4} \text{ s}^{-1}$. The TEM micrographs coincide with the XRD results in Figure 5. A nano-grain structure was observed in all specimens and their grain size was almost constant even after deformation to $\epsilon = 0.28$. In other words, the nano-grain structure shown in Figure 6 was nonuniform during deformation since the compact is in a non-equilibrium state and has the ability to change its structure.

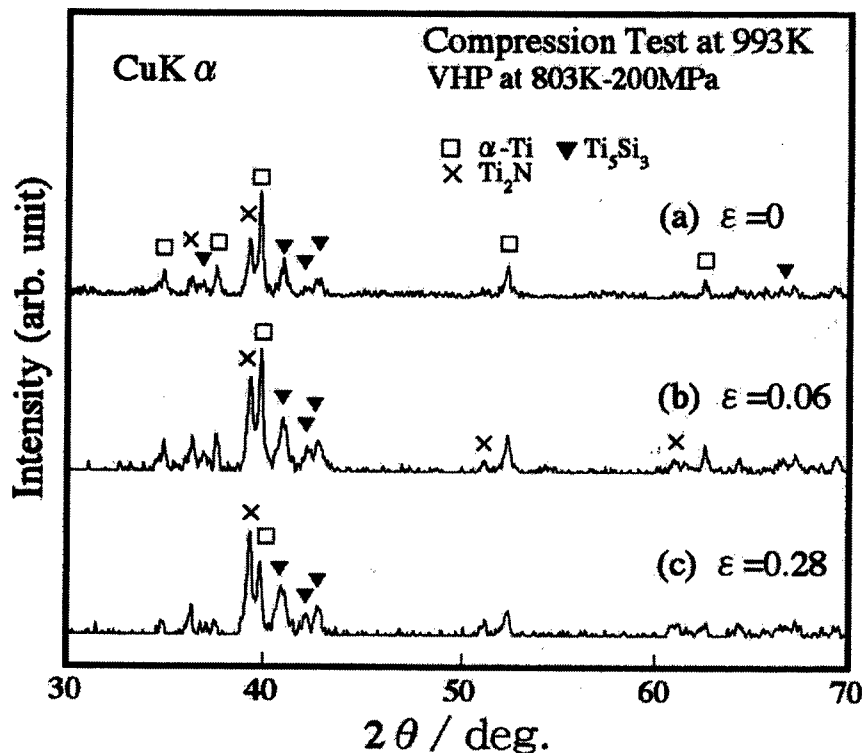


Fig. 5. XRD results of the VHP compact deformed to: a. $\epsilon = 0$; b. $\epsilon = 0.06$; and c. $\epsilon = 0.28$ at 993 K.

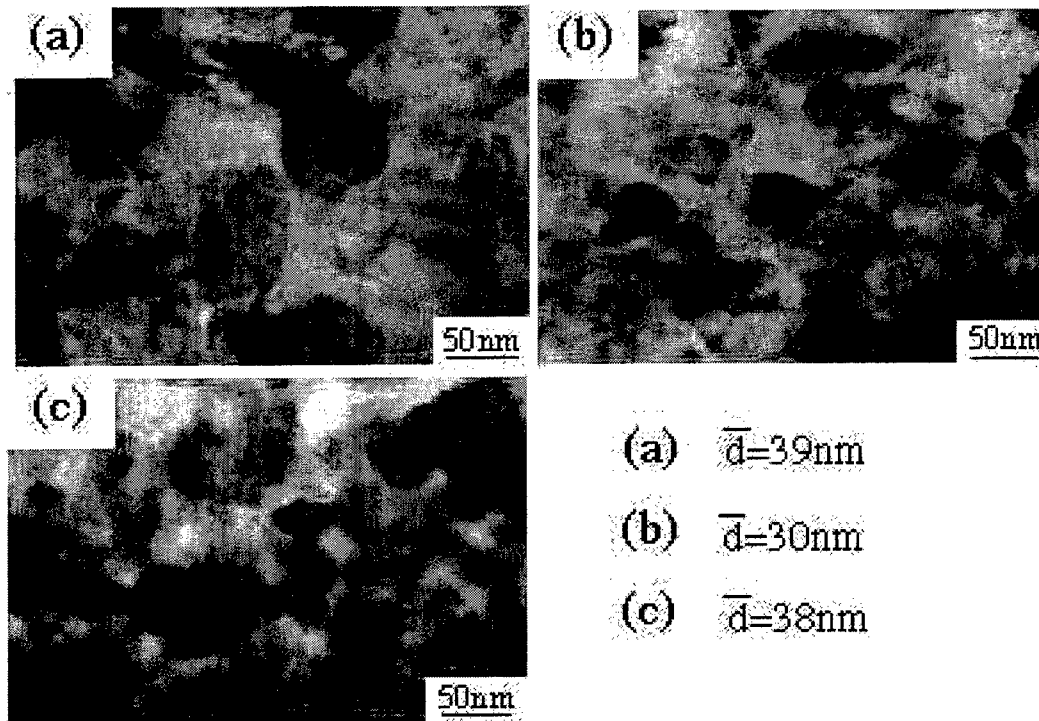


Fig. 6. TEM micrographs of the VHP compact deformed to (a) $\epsilon = 0$, (b) 0.06 and (c) 0.28 at 993 K.

Since the hardness of Ti_2N and Ti_5Si_3 phases are higher than that of the α -Ti phase, increasing these harder phases during deformation resulted in an apparent work-hardening behavior as shown in Fig. 4 (b). The rate of microstructure change in the specimen during high temperature deformation strongly depends on the deformation temperature as well as the strain rate. Therefore, it is expected that there will exist an optimum condition to decrease flow stress for the deformation process.

Figure 7 shows the flow stress of the specimens at $\epsilon = 0.28$ at various temperatures at various initial strain rates. As can be seen, a minimum flow stress was obtained at 993 K at an initial strain rate of $4.2 \times 10^{-4} \text{ s}^{-1}$. Such a reversal in the flow stress between two different initial strain rates was attributed to the microstructural change during deformation. The slower strain rate test produced larger amounts of harder phases such as TiN , Ti_2N and Ti_5Si_3 . In addition, TEM observations demonstrate that the specimen consisted of an (α -Ti + Ti_2N + Ti_5Si_3) microduplex structure with an average grain size of approximately 40 nm.

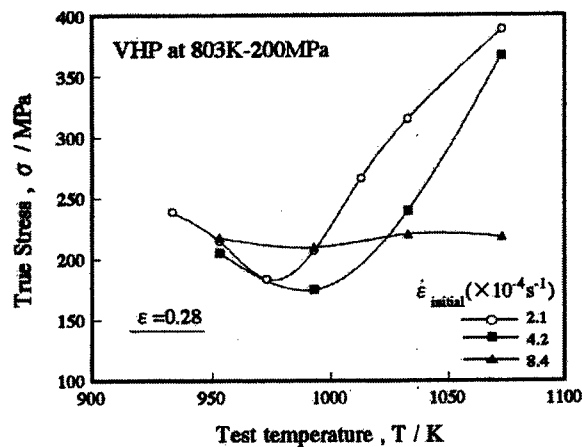


Fig. 7. Flow stress of the specimens at $\epsilon = 0.28$ at various temperatures at various initial strain rates.

Figure 8 shows an XRD result and a TEM micrograph of the specimen deformed to $\epsilon = 0.28$ (i.e., 25% reduction in area) at 993 K at an initial strain rate of $4.2 \times 10^{-4} \text{ s}^{-1}$ followed by heat treatment at 1473 K for 3.6 ks. The XRD and TEM results indicate that the specimen consisted of TiN and Ti_5Si_3 phases and their average grain size was approximately 250 nm. It is well-known that these phases are very hard, so that it requires extremely high temperature for superplastic deformation. However, starting with a non-equilibrium phase compact such as α -Ti and/or an amorphous phase, enables the system to deform at lower temperature and lower flow stress. The same microstructure can be obtained by following with a heat treatment stage. It is noteworthy that the non-equilibrium process enables us to obtain an ultra-fine grain structure. Thus, there exists an appropriate deformation condition to improve the composite forming process as well as mechanical properties of the products.

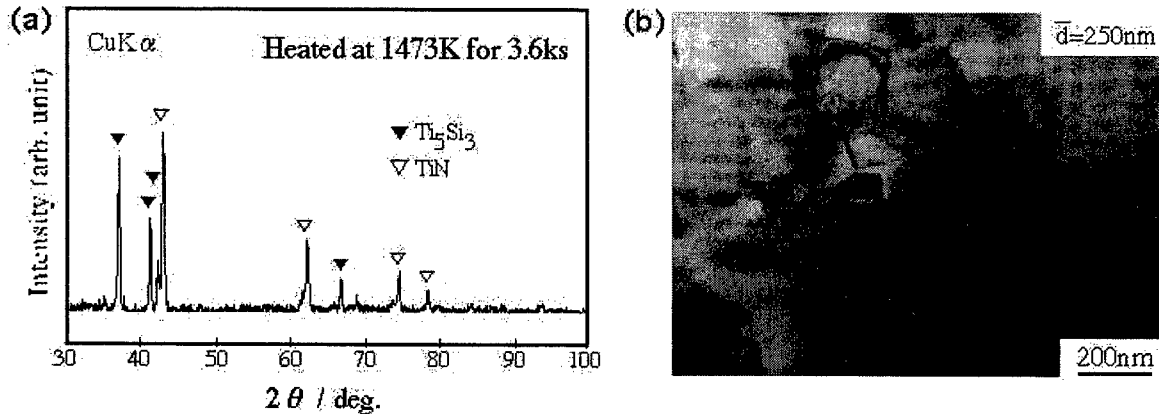


Fig. 8. a. XRD result and b. TEM micrograph of the specimen deformed to $\epsilon = 0.28$ at 993 K at an initial strain rate of $4.2 \times 10^{-4} \text{ s}^{-1}$ followed by heat treatment at 1473 K for 3.6 ks.

CONCLUSION

High temperature deformation behavior and microstructural changes were investigated in a Ti - 20% Si_3N_4 MA powder compact. Deformation of a non-equilibrium phase compact, composed mainly of an α -Ti phase, demonstrated the reversion of the flow stress at elevated temperatures. Such a reverse of flow stress between two different initial strain rates was attributed to a microstructure change during deformation. Utilizing a non-equilibrium deformation process, we are able to improve the composite forming process as well as mechanical properties of the products.

REFERENCES

1. J.S.Benjamin, 1970. Met.Trans., 1, 2943.
2. K.Ameyama, H.Uno, M.Tokizane, 1994. Intermetallics, 2, 315.
3. K.Ameyama, O.Okada, K. Hirai, N.Nakabo, 1995. Mater. Trans. JIM, 36, 269.

Shape Prediction of Growing Billet in Spray Casting Process Using Scanning Gas Atomizer

Eon-Sik Lee*, Woo-Jin Park*, Sangho Ahn* and Shinill Kang**

* Research Institute of Industrial Science and Technology (RIST),
Pohang 790-330, Korea

** Dept. of Mechanical Design and Production Eng., Yonsei University,
Seoul, 120-749, Korea

ABSTRACT

A numerical model has been suggested to predict and analyze the shape of a growing billet produced by spray casting using the scanning gas atomizer. It is important to understand the mechanism of billet growth because a billet with a desired final shape can be obtained by optimum combination of several process parameters. The shape of a growing billet has been determined by the flow rate of alloy melt, spray mass distribution, rotation and withdrawal speed of the substrate, and scanning motion of the gas atomizer. Scanning motion has been controlled by the profile of a cam which determines scanning angle and scanning speed of the gas atomizer. The effects of the most dominant process parameters, such as withdrawal speed of the substrate and the cam profile, on the shape of the growing billet have been discussed. This numerical model can also serve as a basis for heat transfer analysis of the growing billet.

INTRODUCTION

Spray casting is governed by many processing conditions, and therefore is of importance in understanding the effects of such conditions upon the process[1-2]. It is essential to examine the preform growing mechanism because this provides useful information for heat transfer and deformation analyses, and microstructure control, especially when the shape of the preform is three dimensional, e.g. billets. Also, by accurate control of the process, the desired final shape without secondary cutting operations can be produced. The shape of a growing billet is determined by the flow rate of the alloy melt, the mode of gas-atomizer scanning which is due to the cam profile, scanning angle, and the withdrawal speed of the substrate[3-5]. In the present study, a theoretical model was first established to predict the shape of the billet and next the effects of the most dominant processing conditions, such as withdrawal speed of the substrate and the cam profile, on the shape of the growing billet were studied. Process conditions were obtained to produce a billet with uniform diameter and flat top surface, and an ASP30 high speed steel billet was manufactured using the same process conditions established from the simulation.

MATHEMATICAL MODEL FORMULATION

Volumetric Deposition Rate of Spray Cone

It is important to predict mass flux distribution of the spray to calculate the shape of the preform. Fig. 1 shows a schematic of the spray cone. By assuming that all droplets move parallel to the spray axis, the volumetric deposition rate m ($\text{m}^3/\text{m}^2 \text{ s}$) at a point P inside the cone in the direction of spray axis can be written as[6]:

$$m(r_s, d_s) = m_0(d_s) \cdot \exp \left[- \left(\frac{b \cdot r_s}{d_s} \right)^n \right] \quad 1.$$

where d_s is the distance from atomizer to deposition circle, r_s the distance from center of deposition circle to point P, m_0 the volumetric deposition rate at the center of deposition circle which includes the point P, and b is a constant which are determined by experiments. The exponent n is a shape factor of the spray which indicates the extent of the concentration of the droplets on the center of the spray cone.

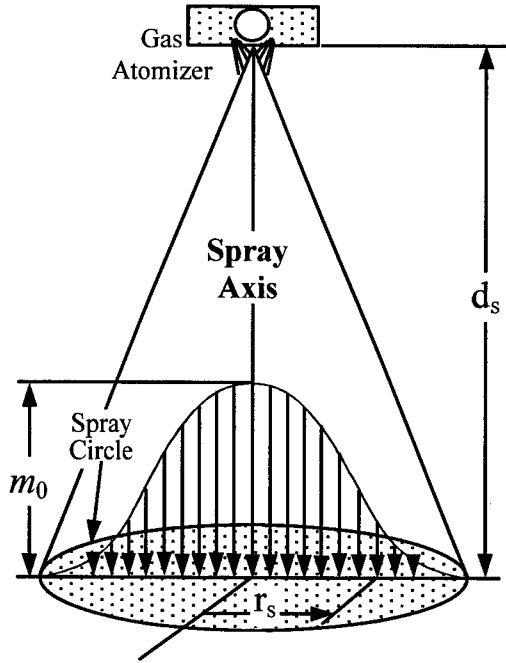


Fig. 1. Schematic droplet mass flux profile in the spray cone

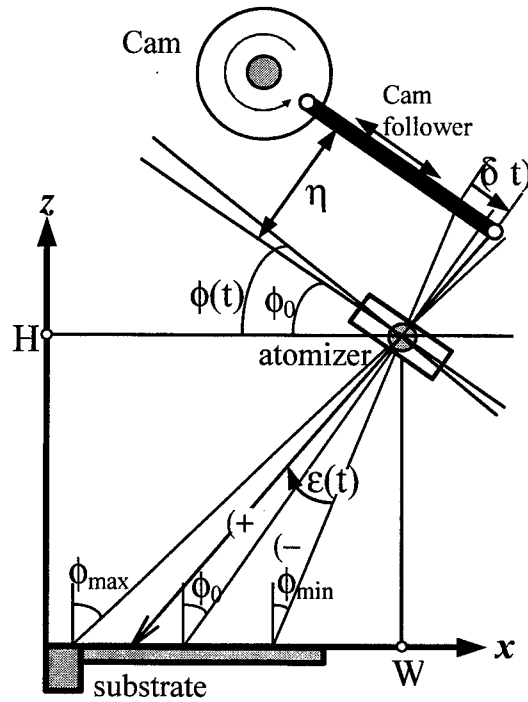


Fig. 2. Scanning motion of the gas atomizer with rotation of disc cam

The volumetric deposition rate at the center of deposition circle can be determined by integrating Eq. (1) as

$$m_0(d_s) = \frac{b^2 \cdot \Phi_v}{\kappa(n) \cdot d_s^2} \quad \text{where} \quad \kappa(n) = \frac{2\pi}{n} \cdot \Gamma\left(\frac{2}{n}\right) \quad 2.$$

and Φ_v means the overall volumetric flow rate (m^3/s) of melt supplied to gas-atomizing nozzle.

SCANNING VELOCITY AND CAM PROFILE

The gas atomizer is scanned by a cam mechanism in order that uniform mass flux and enthalpy of spray across the billet surface are obtained. Fig. 2 shows the scan of gas atomizer with cam rotation. It is necessary to calculate the temporal change of scan angle ϕ of the spray axis as it relates a point on the billet surface with the gas atomizer. The cam displacement δ can be represented in terms of elapsed time t as

$$\delta(t) = \eta \cdot \tan \epsilon(t) \quad 3.$$

where ϵ is the angle between the initial and current spray axes, represented by ϕ_0 and ϕ , respectively. Here, ϕ_0 and ϕ have the relation $\phi(t) = \phi_0 + \epsilon(t)$. The temporal change of scan angle ϕ of the spray axis can be related with δ , ϵ , and η as follows:

$$\tan \phi(t) = \frac{\tan \phi_0 + \tan \epsilon(t)}{1 - \tan \phi_0 \cdot \tan \epsilon(t)} = \frac{\eta \cdot \tan \phi_0 + \delta(t)}{\eta - \tan \phi_0 \cdot \delta(t)} \quad 4.$$

To produce a preform with a desired shape, it is necessary to control the scan of the atomizer. It should be noted that a preform of desired shape can be produced by controlling the scan of the atomizer which affects local deposition rate of droplets. Therefore it is important to design the cam profile, i.e. cam displacement δ , so that a desired scan motion is achieved.

BILLET GROWING MECHANISM

The relative position of the billet, atomizer, substrate, and the position vectors are illustrated in Fig. 3. In a real manufacturing process, the atomizer scans around the billet and the substrate withdraws downwards with speed v while rotating about a vertical axis with angular velocity ω . In modeling the process, we assume that the substrate stands still and that the atomizer rotates about substrate axis and moves upwards while scanning. We also assume the origin of the Cartesian coordinates (X, Y, Z) is attached to the center of the substrate.

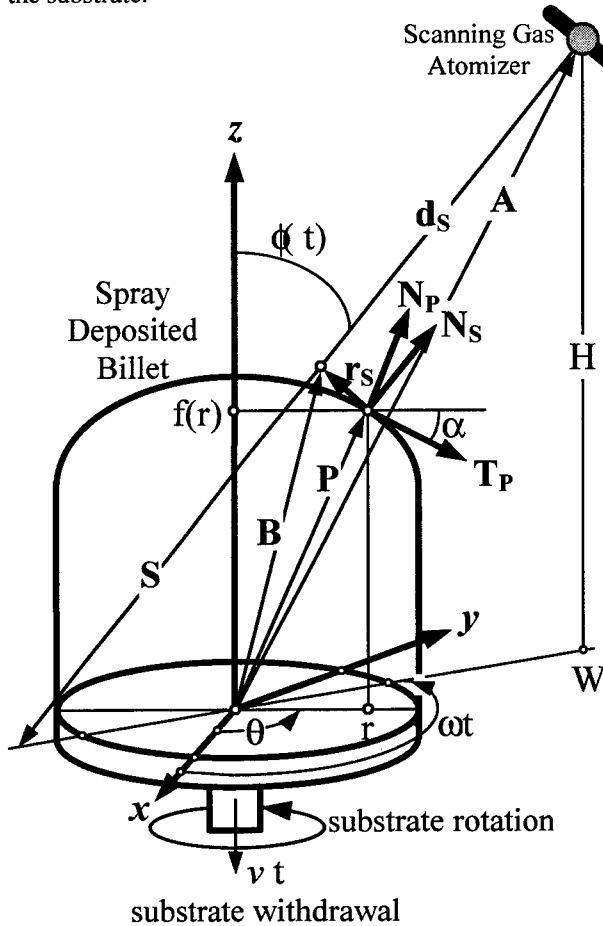


Fig. 3. Position vectors illustrating relationship between gas atomizer, spray cone, growing billet, and moving substrate.

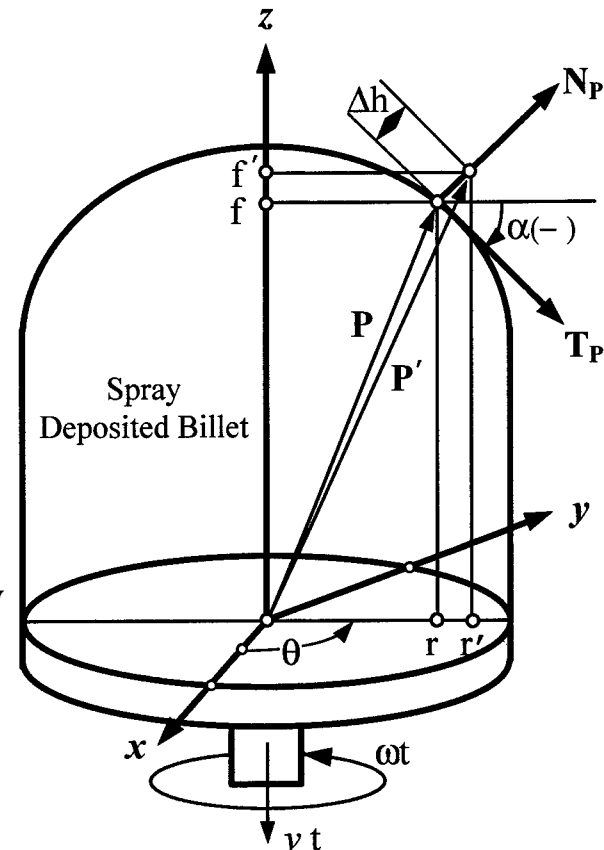


Fig. 4. Incremental change of billet shape during growth of billet.

To predict the deposition rate at each position on the billet surface, several position vectors are defined. From Fig. 3, vectors r_s and d_s can be represented in terms of vectors A, P, S as

$$r_s = B - P = A + \left(\frac{P \cdot S - A \cdot S}{S \cdot S} \right) \cdot S - P \quad 5.$$

$$d_s = B - A = \left(\frac{P \cdot S - A \cdot S}{S \cdot S} \right) \cdot S$$

Vector P indicates arbitrary point P on the billet surface and A the location of atomizer at elapsed time t . Vector S starts from the atomizer and ends at the point at which spray axis intersects X - Y plane. P is independent of time and is expressed in Cartesian coordinates as

$$P = [x, y, z] = [r \cos \theta, r \sin \theta, f(r)] \quad 6.$$

where θ is an angle between X -axis and the plane which includes the point P and Z -axis, r is the radial distance of the point P from Z -axis, and f is the vertical distance of the point P from the substrate. Vector A is given by

$$\mathbf{A} = [W, 0, H + vt] \cdot \begin{bmatrix} \cos \omega t & \sin \omega t & 0 \\ -\sin \omega t & \cos \omega t & 0 \\ 0 & 0 & 1 \end{bmatrix} = [W \cos \omega t, W \sin \omega t, H + vt] \quad 7.$$

where H and W are the initial height of the atomizer and the radial distance of the atomizer from Z -axis, respectively. The vector \mathbf{S} can be represented in terms of the vector \mathbf{A} and the scan angle $\phi(t)$ as

$$\begin{aligned} \mathbf{S} &= [-(H + vt) \tan \phi(t), 0, -(H + vt)] \cdot \begin{bmatrix} \cos \omega t & \sin \omega t & 0 \\ -\sin \omega t & \cos \omega t & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ &= [-(H + vt) \cos \omega t \cdot \tan \phi(t), -(H + vt) \sin \omega t \cdot \tan \phi(t), -(H + vt)] \end{aligned} \quad 8.$$

Now, \mathbf{d}_s and \mathbf{r}_s are obtained by substituting Eqs. (6) to (8) into Eq. (5) as

$$|\mathbf{d}_s| = [(H + vt - f) \cos \phi(t) + [W - r \cos(\omega t - \theta)] \cdot \sin \phi(t)] \quad 9.$$

$$|\mathbf{r}_s| = \sqrt{W^2 + r^2 - 2Wr \cos(\omega t - \theta) + (H + vt - f)^2 - (\sin \phi(t) \cdot [W - r \cos(\omega t - \theta)] + \cos \phi(t) \cdot (H + vt - f))^2} \quad 10.$$

Furthermore, unit normal vector \mathbf{N}_p at point P and unit vector \mathbf{N}_s which is parallel to spray axis are given by

$$\mathbf{N}_s = -\frac{\mathbf{d}_s}{|\mathbf{d}_s|} = [\cos \omega t \cdot \sin \phi(t), \sin \omega t \cdot \sin \phi(t), \cos \phi(t)] \quad 11.$$

$$\mathbf{N}_p = [-\cos \theta \cdot \sin \alpha, -\sin \theta \cdot \sin \alpha, \cos \alpha] \quad 12.$$

where α is an angle between the tangential vector \mathbf{T}_p and the Z -plane, and is defined as $\tan \alpha = df/dr$.

The deposition rate at a point P in the direction normal to billet top surface is calculated as

$$\begin{aligned} \Psi(r, \theta, t) &= [m(r_s, d_s)] \times [\mathbf{N}_s \cdot \mathbf{N}_p] \\ &= \frac{b^2 \Phi_v}{\kappa(n) d_s^2} \cdot \exp \left[-\left(\frac{b r_s}{d_s} \right)^n \right] \times \left\{ \begin{aligned} &\cos \alpha \cdot \cos \phi(t) - \sin \alpha \cdot \cos \theta \cdot \cos \omega t \cdot \sin \phi(t) \\ &-\sin \alpha \cdot \sin \theta \cdot \sin \omega t \cdot \sin \phi(t) \end{aligned} \right\} \end{aligned} \quad 13.$$

The average growth thickness Δh in the direction of normal to the billet surface is obtained by integrating Eq. (13) as

$$\Delta h = \int_t^{t+\Delta t} \Psi(r, \theta, t) dt \quad 14.$$

Fig. 4 shows the incremental change of billet shape. It is noted that position vector P at time $t + \Delta t$ is updated from an arbitrary billet surface at time t as follows

$$\begin{aligned} \mathbf{P}' &= [x', y', z'] = \mathbf{P} + \Delta h \cdot \mathbf{N}_p \\ &= [(r - \Delta h \cdot \sin \alpha) \cos \theta, (r - \Delta h \cdot \sin \alpha) \sin \theta, f(r) + \Delta h \cdot \cos \alpha] \end{aligned} \quad 15.$$

RESULTS AND DISCUSSION

Table 1 shows the process conditions established from a series of computer simulations to produce a billet which has uniform diameter of 0.14 m and flat top surface while maintaining a constant growth rate. The corresponding billet shape calculated using the process conditions in Table 1 is shown in Fig. 5. This indicates the possibility of producing an actual billet with uniform morphology and scale of microstructures, by controlling the local growth rate at the billet top surface.

Effect of Substrate Withdrawal Speed

Process simulations have been performed to observe the effect of substrate withdrawal speed on the billet shape. Billet shapes were calculated using the same process conditions as in Table 1 except the substrate withdrawal speed. It was changed to $0.00025 \text{ m} \cdot \text{s}^{-1}$ (case A) and $0.00055 \text{ m} \cdot \text{s}^{-1}$ (case B), and the

corresponding billet shapes are shown in Fig. 6a (case A) and Fig. 6b (case B). In case that the withdrawal speed is too slow compared with the billet growth rate as in case A, spray axis tends towards edge of the billet surface. As a result, the shape of the top surface becomes concave and the billet diameter becomes larger as deposition process continues. In contrast to case A, if the withdrawal speed becomes too fast compared with the billet growth rate, the billet diameter gets smaller and the shape of the top surface becomes convex, as in case B. Therefore, to manufacture a billet with uniform diameter and flat top surface, one should maintain the distance between atomizer and deposition surface constant by controlling the process so that the substrate withdrawal speed and the billet growth rate are balanced.

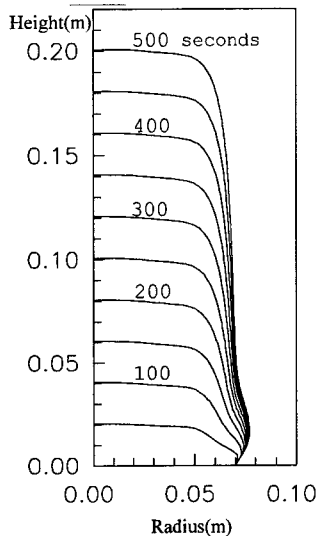


Fig. 5. Desired billet shape at each elapsed time.

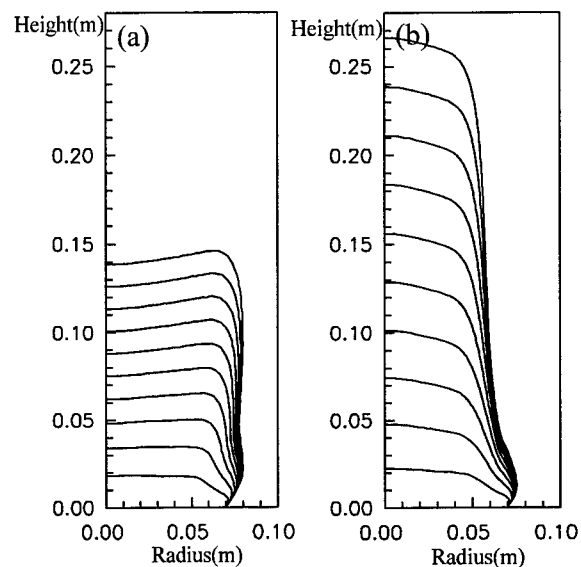


Fig. 6. Billet shape variation due to different withdrawal speeds.(a) $v=0.00025\text{m/s}$ (b) $v=0.00055\text{m/s}$

Effect of Scanning Cam Profile

In the present research, a scanning cam mechanism was designed to regulate density and enthalpy of the spray. Fig. 10 is the plot of the radial velocities and displacements of the scanning cams used for the present simulations to study the effect of cam profile on the billet shape. Radial velocity of cam 1 was designed so that the gas atomizer scans slower at the edge than at the center of the billet top surface. As a consequence, a billet with flat top surface could be obtained throughout the whole deposition process as already shown in Fig. 5. However, in case of cam 2, the atomizer has the same scan velocity at the center and at the edge, and it results in higher growth rate at the center than the edge of the billet top surface. Fig. 8 shows the billet shape obtained using the radial velocity of cam 2. Due to higher growth rate at the center, billet top surface is very uneven and the radius of the billet becomes smaller as the process continues. However, after 150 rotations high growth rate at the center exceeds the substrate withdrawal speed, and it results in increase of billet radius. It is also observed that the local deposition rate varies as deposition proceeds, and it may yield a billet of non-uniform microstructures.

CONCLUSIONS

In the present study, a numerical method was presented to predict and analyze the shape of a growing billet produced from the spray forming process. A theoretical model was first established to predict the shape of the billet and then the effects of the most dominant process conditions, such as withdrawal speed of the substrate and the cam profile, on the shape of the growing billet were studied. The design guidelines were extracted from the simulation to manufacture a billet with uniform diameter and flat top surface as follows:

- (1) The process should be controlled so the substrate withdrawal speed and the billet growth rate are balanced.
- (2) The scanning cam mechanism should be designed so the gas atomizer scans slower at the edge than at the center of the billet top surface.

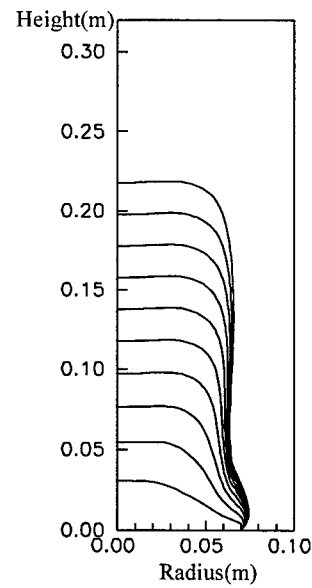
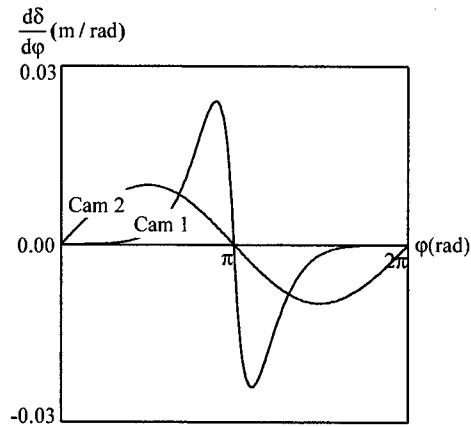


Fig. 7. Cam radial velocities for Cam 1 and Cam 2. **Fig. 8.** Billet shape produced using Cam 2.

Table 1. Process conditions to produce desired billet shape

$\Phi_v = 4 \times 10^{-5} m^3 s^{-1}$	$n = 1.4$	half scan angle = 3°
$H = 0.25 m$	$W = 0.20 m$	$\phi_0 = 34^\circ$
$\omega = 2\pi rad s^{-1}$	$\omega_c = 2\pi \times 8.1 rad s^{-1}$	$v = 0.0004 m s^{-1}$

REFERENCES

1. Eon-Sik Lee, W.J. Park, K. H. Baik, S. Ahn, 1998. Scripta Materialia, 39(8), 1133 ~ 1138.
2. Eon-Sik Lee, W. J. Park, J.Y. Jung, S. Ahn, 1998. Met. Mater. Transactions A, 29(5), 1395-1404.
3. L. Warner, C. Cai, S. Annavarapu, R. Doherty, 1997. Powder Metallurgy, 40, 121-125.
4. Eon-Sik Lee, S. Ahn, Shiill Kang, 1997. J. of the Korean Inst. of Met. & Mater., 35, 460-467.
5. Shinill Kang and D.-H. Chang, 1999. Materials Sci. and Eng., A260, 161-169.
6. J. Forest, S. Lile, J. Coombs, 1993. Proc. 2nd Inter. Conf. Spray Forming, 117-127.

Intelligence in Concurrent Engineering

Modelling Design Planning in Concurrent Engineering

C. Reidsema and E. Szczerbicki

Department of Mechanical Engineering,
University of Newcastle,
Callaghan 2308, Australia

ABSTRACT

Concurrent Engineering offers to industry, the promise of significant reductions in product development times matched to equally significant improvements in product and process quality. To adequately realise this promise, we suggest that a crucial focal point is the development of models that further the understanding of the nature of design planning within such complex environments. Progress in the long term towards developing computer-based techniques and tools to support human participants within the organisation may only proceed with this understanding. Although computer based systems are increasingly being utilised to support process improvements in CE manufacturing planning, there exists a significant gap in the knowledge of not only how this technology can be used to support upstream design planning decisions, but in the fundamental paradigms that would be necessary for the development of such a design process planning system. In this paper we suggest that the Blackboard Database is an appropriate tool for investigating and developing models and strategies to represent human cognition in CE design planning.

INTRODUCTION

Concurrent Engineering (CE) as a philosophy for product development, emphasises two central themes: namely, (a) simultaneous execution of tasks in order to minimise development time; and (b) utilisation of downstream life-cycle knowledge in upstream development processes to maximise product and process quality. Although satisfaction of either objective may decrease the project cycle time, it is generally not known *a priori* whether selection of one strategy over another will produce the most favourable result with respect to other constraints and objectives existent at different levels of an organisation. For example, a strategic objective of an organisation may be to secure market share through quality and the company may be willing to absorb project cost overruns or extensions to achieve it. The authors suggest that emphasis should not be simply to either maximise task concurrency or maximise utilisation of life-cycle knowledge. Rather, emphasis should be on maximising the effectiveness of resource utilisation so that design tasks execute at the right time, for the right reason, meeting the right requirements and giving the right results [1].

Within the product development process, design presents itself as perhaps the most complex, least understood, yet most promising area for leveraging significant gains in terms of cost savings, risk reduction, and quality improvements through planning [2,3]. Planning is essential to successful leverage of these CE benefits because of the dynamic nature of the organisation, the product, and the design process itself. Planning can be viewed as a decision-making process that attempts to predetermine a course of action to achieve some particular goal [4]. Although product design often consists of intuitive and creative-thought not amenable to computerisation, we take the broader view that design and design planning are primarily decision-making processes that are largely repetitive tasks which contain identifiable reasoning structures, strategies, constraints, rules, and data that can be automated. These basic decision-making elements form the basis of models that represent design planning activities of a human planner in a CE design environment.

An artificial intelligence tool called a blackboard database provides a way to examine and formulate coherence of the basic planning elements by viewing them as a system model containing: (a) an abstracted problem state, (b) knowledge that acts upon the problem state, and (c) strategies to apply the knowledge to the problem state. We suggest that an appropriate model to capture essential features of human-centric

planning in a distributed CE design environment is one that represents the problem state in terms of four primary levels of abstraction referred to as GDDI [5]:

- (1) Plan generation (G),
- (2) Plan decomposition (D),
- (3) Plan distribution (D), and
- (4) Plan Integration (I).

Knowledge sources consist of production rules and algorithms applied to solve the individual sub-problems at each level. A solution for one sub-problem becomes the input to another sub-problem and so on, until a complete solution is obtained. This problem-solving activity is monitored and coordinated by the control source which may be encoded with knowledge application strategies or reasoning. We present additional models which may be used to devise knowledge sources and control strategies that can be applied to provide solutions for each sub-problem contained within the abstracted GDDI problem state.

MODELLING CE DESIGN ENVIRONMENT

The design phase appears to be the most predominant aspect of product development in that its influence on the cost and functionality of a finished product and thus an organisation's profitability, is grossly disproportionate to the investment cost [6]. Traditional methods for structuring resources within an organisation to design, develop, manufacture and distribute products to customers are proving inadequate in an economic climate exemplified by diminishing windows of opportunity for time-to-market and demand for ever-increasing quality improvements. Departmentalisation and limited cooperation is giving way to a reorganising of expertise and functions into distributed regimes of experts who may not even reside within the same national boundaries of the parent company. These new organisational structures are product-focused in that the needs of the product at any point in its development dictate the resource structure of the organisation. This holistic view is partially achieved through use of Product Development Teams (PDTs). A PDT may be thought of as a group of individuals or agents whose collective knowledge and expertise can be efficiently matched to a defineable set of problems representing a specific phase of the product development process [7]. The planning of these team-activities is highly complex requiring coordination of inputs to decision-making from many areas within the organisation which must be rationalised to prevent an overflow of contributions from slowing down the design process [8].

The knowledge available to PDTs to support design planning can be modelled as either Product, Process or Organisational (PPO) knowledge which may be in the form of requirements and constraints formulated as data and rule sets which assist in guiding decision-making [9]. The types of knowledge that is found within each category of the knowledge model are numerous, varied and highly-dependent on the purpose to which it is utilised. Figure 1 shows a few of the knowledge types that pertain specifically to design planning.

Product Knowledge	Process Knowledge	Organisational Knowledge
Life-cycle Attributes	Requirements	Strategic
Product Specification	Conceptual	Tactical
Product Hierarchical Structures	Analysis	Operational
	Embodiment	

Fig. 1. Product, Process and Organisational Knowledge Types for Design Process Planning

The Process Knowledge shown in Figure 1 consists of a number of distinct decision-making phases that represent a common prescriptive model of the design process [10]. In a logical sense, design phases are useful to represent transformations of knowledge into a product and so, these reasoning are used to determine the state of the product along its path of evolution. These heuristics provide a means to determine the knowledge needed, when, and what decisions have been made and which decisions are yet to be made.

The design process does not exist in isolation. Rather it operates within an organisation with its own primary imperatives and goals concerned with maintaining competitive advantage. The Organisational Knowledge shown in Figure 1 depicts three well-established levels of knowledge within a hierarchical authority structure that bears significantly on design process planning. For example, a company may have a competitive strategy with the objective to develop a distinctive quality advantage for its product. A tactic to accomplish this goal may be to specify that a reliability analysis of several key components of the product be conducted. At the operational level this may entail generating additional tasks in a particular sequence as well as allocating resources to these tasks for their completion. Determining which tasks must be performed and the degree to which these tasks are accomplished without expensive project overruns is at the core of CE design planning.

PLANNING MODELS

There are a number of reasonable definitions for planning available in the literature [4,11,12]. A common denominator in all these definitions, is that planning is a cognitive decision-making process to specify a course of action to achieve a specified objective [5]. The term "cognitive" is crucial to our research in that it emphasises that planning is viewed as a problem in knowledge application and the determination of strategies to apply this knowledge should be our primary goal. One strategy which appears to emulate satisfactorily the cognitive problem-solving behaviour of human planners is referred to as "opportunistic strategy" [3,4,7]. In such a strategy, decisions on what knowledge and when it should be applied is based on the current state of the solution and whether or not it will be advanced through this action. An opportunistic strategy is well suited to complex, ill-structured, indeterminate problems such as design planning.

We can model decision-making in design planning by expanding the decision-making model of Intelligence, Design and Choice proposed by Simon [13]. In this model we:

1. define an objective from, for example, specifications and requirements;
2. identify or search for tasks which when temporally ordered will satisfy the objective;
3. generate a number of these plan sets;
4. test these plan sets for effectiveness in satisfying the objective by using known constraints; and
5. make a choice as to which plan set is preferable.

The salient points that we derive from this simple expansion are concerned with selecting appropriate models for our objectives and the required tasks. There are several possible models available such as Program Evaluation and Review Technique (PERT), Structured Analysis and Design Technique (SADT or IDEF0), and Design Structure Matrices (DSM). Of these techniques we feel that SADT provides the most promising starting point to model tasks for the purpose of design planning. At the core of SADT is the Structured Analysis (SA) box shown in Figure 2 [6].

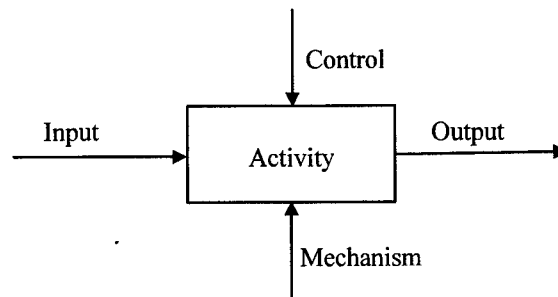


Fig. 2. Structured Analysis (SA) Box

The SA box represents an activity which transforms inputs into desired outputs and is guided and constrained by control data. Mechanisms are those elements which are required to support the activity but more importantly may be used to cross-reference models (organisational, resource etc...) which are needed for task completion. While all of these techniques have advantages in representing tasks in design planning,

further investigation will be necessary to adapt the SADT task model to satisfy the structural requirements imposed by the blackboard system.

BLACKBOARD REPRESENTATION

The Blackboard Database Architecture (BBDA) is a knowledge-based problem solving system which emulates the real-life scenario of experts who, in a controlled manner, cooperate by contributing partial-solutions to a larger complex problem represented on a blackboard [14]. The BBDA has been extensively utilised as a tool for investigating and representing complex, distributed problems with no determinate solution such as those found in concurrent design process planning [15,16]. The three main components of a BBDA which are required to represent this metaphor are: the Blackboard Database, Knowledge Sources (KSs) and a Control Source (CS).

The Blackboard Database acts as centralised storage for data and information placed there by KSs. These KSs may interact with each other only through the database which also provides a common data structure for integrated communication. The database is hierarchically structured into separate, related levels of abstraction, each corresponding to partial problems that are loosely-coupled to form the overall problem. The knowledge required to solve these partial problems are contained within KSs in the form of production rules or algorithms designed specifically for each partial problem. The CS makes a number of decisions on KS coordination, including, but not limited to: (1) which KS will be selected out of those that are triggered, (2) when this KS will be executed, and (3) where in the problem structure it will be applied. A CS may operate opportunistically or be encoded with predetermined knowledge application strategies to enforce or induce a preferred system behaviour [5]. It is primarily this flexibility of problem-solving behaviour that makes the BBDA a promising candidate for modelling the cognitive behaviour of human planners.

Determining the structure of the blackboard database requires the selection of an abstraction hierarchy that will best represent the overall problem of cooperative planning in a distributed environment from the most global perspective [17]. The four primary levels of abstraction of GDDI referred to in the introduction make up the top-level of the blackboard database hierarchical structure as shown in Figure 3. We depict data received from external agents such as PDT's, management and other participants in the planning process serving as input to the top-level Integration abstraction layer of the hierarchy. An Integration KS transforms data on this level to provide input to the Generation level which then becomes input to the Decomposition level and so on, until a fully distributed plan consisting of partial plan sets with associated constraints are transmitted as output data to the external agents in the system. These agents may execute these plans and modify constraints and the cycle will repeat itself ensuring that a current global plan is maintained.

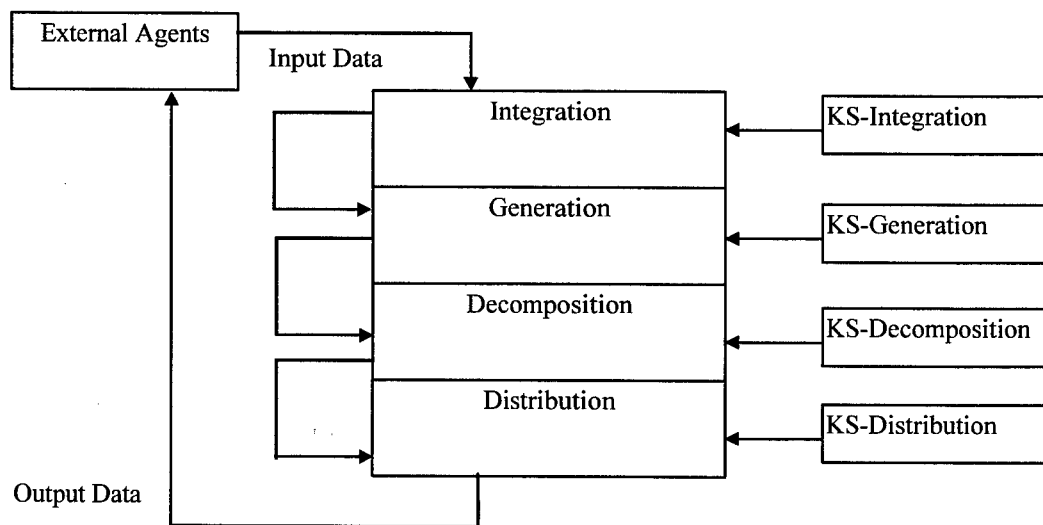


Fig. 3. Blackboard Database Structure for Planning

FURTHER RESEARCH

Our research is currently focused on the design and specification of the planning problem as it will be represented within the blackboard database data structure. This entails asking ourselves a number of questions, such as:

1. What constitutes an acceptable solution to the problem?
2. In what form should the solution take as output of the BBDA?
3. What information do we need as input to each level of the GDDI solution? and
4. What knowledge is needed to produce output for each level of the GDDI solution?

We anticipate that answers to these questions will best be facilitated by designing a simplified design planning problem model and interpreting the information that is contained within the data of this model. We expect that a rigorous examination of this model will provide us with a clearer understanding of the knowledge required to:

- a. generate plans;
- b. decompose these plans into sets of tasks;
- c. distribute these tasks to agents for execution and modification; and
- d. retrieve and integrate the results from these agents to generate a new plan.

A primary expectation of this research is the development of a simplified working blackboard system to prove the benefits of providing engineering management with a design process planning software decision support tool that is based on a clear understanding of the CE design planning process in a distributed product development environment.

REFERENCES

1. A.H.B. Duffy, 1995. Ensuring Competitive Advantage with Design Coordination.. 2nd International Conf. on Design to Manufacture in Modern Industry Proceedings - Part 1, Slovenia - Bled.
2. A.H.B. Duffy, M.M. Andreasen, K.J. MacCallum and L.N. Reijers, 1993. Design Coordination for Concurrent Engineering. *Journal of Engineering Design*, 4(4), 251-265.
3. C. Reidsema, E. Szczerbicki, 1998. Blackboard Approach in Design Planning for Concurrent Engineering Environment. *Cybernetics and Systems: An International Journal*, 29(7), 729-750.
4. B. Hayes-Roth, F. Hayes-Roth, 1979. A Cognitive Model of Planning. *Readings in Planning*, James Allen, James Hendler, and Austin Tate (Eds.), Morgan Kaufmann Pub., Inc. San Mateo, CA, 245-262.
5. C. Reidsema, E. Szczerbicki, 1998. Blackboard Approach for Concurrent Engineering Design Process Planning. *Proc. 1998 Pacific Conference on Manufacturing*, August 18-20, Brisbane, Queensland, Australia. 522-527.
6. S.D. Eppinger, D.E. Whitney, D.A. Gebala, 1992. Organizing the Tasks in Complex Design Projects: Development of Tools to Represent Design Procedures, NSF Design and Manufacturing Systems Conf., Atlanta, Georgia. 301-309.
7. K.J. MacCallum and I.M. Carter, 1992. Opportunistic software architecture for support of design coordination. *Knowledge-Based Systems*, 5(1), 55-65.
8. R.P. Smith, S.D. Eppinger, 1997. A Predictive Model of Sequential Iteration in Engineering Design. *Management Science*, 43(8), 1104-1120.
9. B. Prasad, 1996. *Concurrent Engineering Fundamentals: Integrated Product Development*, II, Prentice-Hall.
10. N Cross, 1989. *Engineering Design Methods*, John Wiley & Sons.
11. R. Tebay, J. Atherton, S.H. Wearne, 1984. Mechanical engineering design decisions: instances of practice compared with theory, *Proc Instn. Mech. Engrs.*, 198B(6).
12. B. Hayes-Roth, F. Hayes-Roth, F. Rosenschein, S. Cammarata, 1979. Modelling Planning as an Incremental, Opportunistic Process. *Proc. 6th International Joint Conference on Artificial Intelligence*. Los Altos, California, William Kaufmann, Inc., 375-383
13. R.H. Sprague, Jr., 1980. A Framework for the Development of Decision Support Systems, *MIS Quarterly*, 4(4), 7-31.

14. H.P. Nii, 1986. Blackboard Systems: The Blackboard Model of Problem Solving and the Evolution of Blackboard Architectures, The AI Magazine, Summer, 38-106
15. N.M. Sadeh, D. W. Hildum, T. J. Laliberty, J. McANulty, D. Kjenstad and A. Tseng, 1998. A Blackboard Architecture for Integrating Process Planning and Production Scheduling. Concurrent Engineering: Research and Application. 6(2).
16. J-P. Kruth, G. Van Zeir, J. Detand, 1996, An Interactive CAPP Kernel Based on a Blackboard System Architecture. Proceedings of The 1996 ASME Design Engineering Technical Conferences and Computers in Engineering Conference, Irvine, California.1-13.
17. I.D. Craig, 1995. Formal Techniques in the Development of Blackboard Systems. Research Report No. CS-RR-199. Department of Computer Science. University of Warwick. Coventry. UK.,
<http://www.dcs.warwick.ac.uk/pub/reports/rr/199.html>

Computer-Aided Integrated Design for Injection Molding

Y.-M. Chen *, C. T. Ho **, R.-S. Lee ***

* Institute of Manufacturing Engineering,
National Cheng Kung University, Tainan, Taiwan, ROC

** Department of Industrial Engineering and Management,
National Kaohsiung Institute of Technology, Kaohsiung, Taiwan, ROC

*** Department of Mechanical Engineering,
National Cheng Kung University, Tainan, Taiwan, ROC

ABSTRACT

A high quality, cost effective product can not be consistently and economically produced unless product design is developed based on a comprehensive understanding of all relevant factors and their interactions. This research aims to: (1) develop an integrated design for injection molding methodology that accommodates concepts of concurrent engineering; and (2) develop a computer-integrated design for injection molding based on this methodology to provide interactive design aids and consistent design assessment to improve injection molding product design quality, producibility, and cost effectiveness.

INTRODUCTION

Injection molding is a significant net shape manufacturing process for producing high tolerance, precisely defined plastic components by forcing molten plastic under high pressure into a split material mold. Injection molding product design may start with a new design or a redesign for injection molding from a preliminary design. They all involve a wide variety of design expertise, knowledge of engineering and injection molding processes and the use of computer-aided tools [1]. As many details need to be considered simultaneously, a great deal of trial and error occurs in the design process. Moreover, due to the complexity, details are often overlooked even by experienced designers. As a result, the design quality is not consistent. And, since redesign or design modifications are made frequently to make a design manufacturable, injection molding product design is a long and costly iterative process.

In recent years, *concurrent engineering* has emerged to improve product competitiveness by resolving product life cycle concerns at the earliest stages of design. These concerns interrelate the entire product life cycle from conceptual design through manufacturing to disposal, including product functionality, producibility, assembleability, serviceability, and even, recycleability. One of the practices of concurrent engineering is the use of computer-aided design for X tools for functionally acceptable, manufacturable, assembleable, and recycleable design [2, 3, 4]. Most of the methodologies are analogous to the method of design for manufacturing (DFM) and design for assembly (DFA) [5] which concentrated merely on a specific life cycle concern, such as product producibility, functionality, or cost effectiveness.

Instead of focusing on a particular development concern, the design for injection molding research being conducted at the Computer-Aided Concurrent Engineering Research Lab. of National Cheng Kung University addresses molding product life cycle issues, including product functionality, quality, producibility, tooling, and cost, in an integrated fashion at each product design stage.

CHARACTERIZATION OF DESIGN FOR INJECTION MOLDING

The molding product life cycle includes the activities of product design, process design, mold design, and mold manufacturing process planning. The molding product and process development concerns that should be considered in the product design stage include: product requirements and specifications, producibility, quality, time to market and cost. Each of these concerns has interdependencies with other concerns.

Points considered in product requirements and specifications include product functions, mechanical and

physical properties, dimensions and tolerances, surface finish, precision tightness, hardness, etc. Principles related to product requirements and specifications include:

1. avoid unnecessary tight/high specifications;
2. relate critical dimensions to only one mold member, and
3. minimize draft angles in critical material content areas.

Product lead-time is determined by product and process development time, tool fabrication time, and product manufacturing time. Similarly, the cost of a molded product primarily consists of material costs, processing costs, post-processing costs, and tooling costs. The factors affecting product time to market and cost are inter-related. A product configuration optimized for injection molding should:

1. fill completely with material;
2. solidify quickly and without defect;
3. eject readily from the mold; and
4. release heat smoothly.

FRAMEWORK OF INTEGRATED DESIGN FOR INJECTION MOLDING

The integrated design for injection molding framework was developed based on the approach of feature-based design, the principles of concurrent design and the methodologies of interactive design evaluation, iterative redesign, and life cycle concern assessment.

In *feature-based design*, a product is constructed, edited, and manipulated in terms of features with certain spatial and functional relationships. At the conceptual design stage, product requirements are specified in terms of functions, each of which corresponds to one or more physical components called form features or functional features. In preliminary design, the major shape of a product is configured using the functional features based on the products functional requirements. The features themselves are functionally defined by attributes which represent significant design aspect and manufacturing semantics, and which are geometrically represented by a set of parameters [6, 7, 8].

The basic idea of *concurrent design* is that product life cycle issues are considered and reviewed throughout all product and process development stages. In the proposed framework, the concept of concurrent engineering is implemented by (1) providing design advice at each product design step based on evaluations against molding product design concerns to avoid unacceptable designs and (2) performing life cycle assessments for mold design and process design on product design results to ensure that all life cycle issues are resolved at the product design stage.

Instead of fully automating the design and redesign, an *interactive design evaluation and iterative redesign* approach is employed. That is, the product design is performed as an iterative "design, design evaluation, and redesign" process. The product designer is responsible for initial design decisions and the design advisory system is interactively evaluating each decision. If any violation occurs, a redesign is suggested and the designer is responsible for the redesign. The iteration continues until no further violations exist or the designer accepts the result.

Based on the methodologies, a framework for integrated molding product design is proposed as shown in figure 1. It includes the phases of molding product design, molding process design concern assessment, molding mold design concern assessment, and mold fabrication concern assessment. In the phase of product design, concerns of product functionality, producibility, quality, and cost effectiveness are evaluated at each design stage. The product design result is then assessed against the concerns of process design, mold design and mold fabrication.

The activities of molding product design include feature design, preliminary design, parting line design, and detail design. The objective of feature design is to design a feature that is functionally acceptable, as well as producible and cost effective. At this stage, the product designer is responsible for selecting a feature that fulfills product functionality. Evaluations for producibility and cost effectiveness of the part feature and manufacturability of corresponding mold feature(s) are performed. The feature placement is

then conducted, which is followed by an evaluation of the interactions between the placed feature and other interacting features. The evaluation focuses on the impact caused by the interactions on product producibility, cost effectiveness, and mold manufacturability. After creation of initial product geometry, a global shape evaluation is performed to ensure producibility of the product and mold manufacturability.

Parting line design includes parting line specifications and evaluations of product producibility and cost effectiveness and mold manufacturability according to the specifications. Undercut detection and mold cavity manufacturability assessments are the major tasks in this step.

Detailed design refines the preliminary product geometry into a moldable final product model by adding drafts and rounding corners based on the specified parting line locations. Evaluations of product producibility, cost effectiveness and mold manufacturability are also conducted at this stage. Life cycle assessments are performed at the end of product design to ensure the design results do not conflict with the concerns of other life cycle activities.

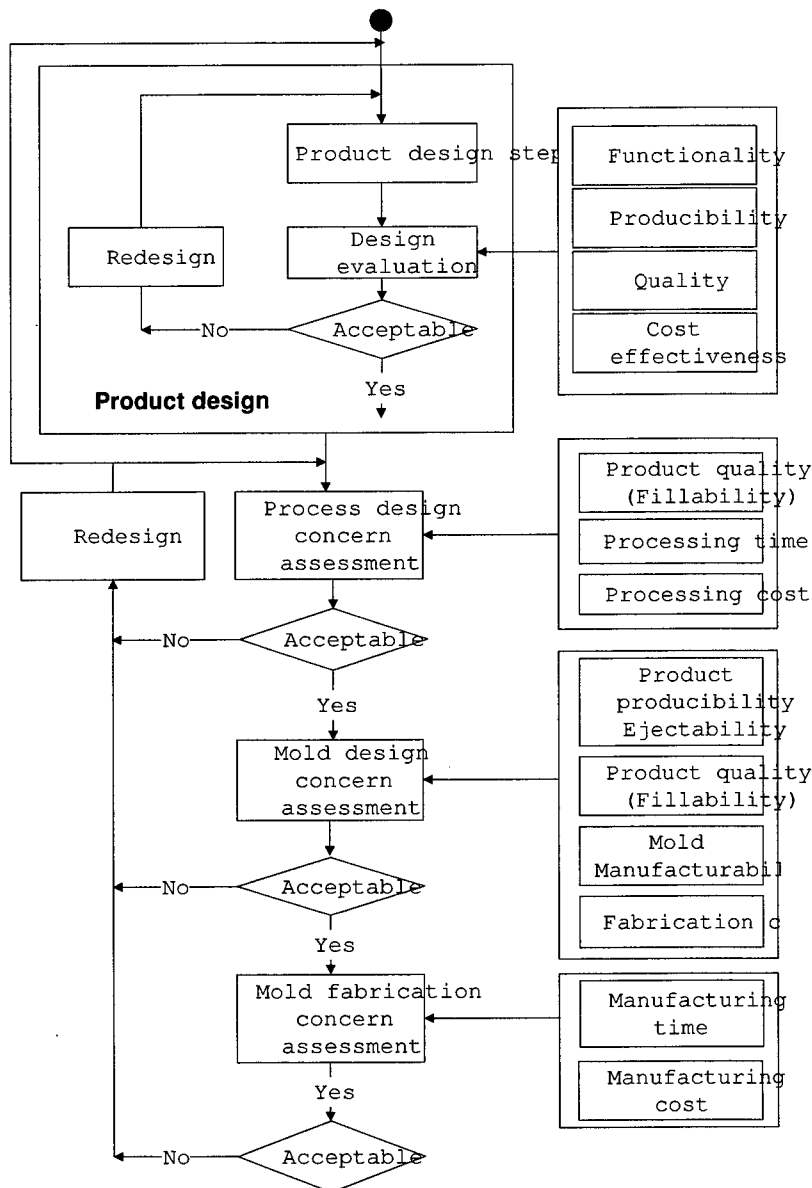


Figure 1. Integrated molding product design framework.

SHAPE CHARACTERISTICS AND PRODUCT DESIGN CONCERN ANALYSIS

Most of the evaluations discussed in the previous section are performed to check if any geometric characteristics or significant items may conflict with design concerns. Therefore, analysis on part geometric characteristics and the relationships between design concerns is one of the major tasks in this research. This section presents the results of the analysis.

Considered from a functional perspective, features can be classified into positive and negative elements [8]. The shape of a feature can be an individual primitive or the combination of several primitives, which can be a solid or a surface. A solid primitive is created in terms of a 2D profile and a trajectory operation. The shape characteristics of features can be categorized based on their operations and trajectories, and the material fillability and relative processing time and cost of each category can be obtained.

A positive part feature forms a cavity in a mold, while a negative feature makes a protrusion in a mold, that can be formed using an insert or core. Mold features and their shape characteristics can therefore be obtained from part features. Similar to part feature characteristics, the shape of mold cavities can be categorized according to their operation and trajectory. Based on trajectory type, a suitable manufacturing method for each category can be identified and prioritized. The relative manufacturing time and cost can be obtained based on the manufacturing methods and number of dimensions for manufacturing operation.

Besides individual feature characteristics, most shape characteristics or significant items are formed by feature interactions that depend highly on spatial relationships between features. Therefore, Abstraction of shape characteristics formed by feature interactions can be done by clarifying spatial relationships between features. The spatial relationships between 2 coplanar features can be classified into "adjacent" or "non-adjacent" relationships. The former indicates 2 features that share at least one face, while the latter indicates two features that share no faces. The "adjacent" relationships can be further classified based on the type of features. Two positive features can be either "adjacent to" or "intersect" each other; and a positive feature can be an "add on" to a negative feature. On the contrary, perhaps a negative feature "is in" a positive feature. Similarly, two negative features can be "adjacent to" or "intersect" each other. To summarize, relationships between features include "adjacent to", "intersect", "add on", "is in", and "coplanar".

Based on the feature spatial relationships, the geometric characteristics can be generalized as significant items, such as depth, thickness, height, volume, and cross-sectional-area. Two adjacent positive features may have a stacked-up "height" or "thickness". The "area", "height" and "thickness" are significant items for two intersecting features. Similarly, "distance" is significant for two coplanar but non-adjacent features.

Adding a positive feature onto a negative feature forms a protrusion with a significant height or thickness. Placing a negative feature into a positive one forms 3 types of cavities - "holes", "grooves" or "steps" depending on the degree of freedom for placement. Each type of cavity owns significant items such as thickness, depth, and cross-sectional area. A "hole" type negative feature can only be placed in one direction to form a "hole" cavity in a part. A "groove" type negative feature can be placed in two directions to form a "slot" or "pocket" cavity. A "step" type negative feature can be placed in three directions and forms a "step" cavity. Placement direction is the orientation that a core can be inserted into and removed from a part or the direction that a cutting tool can approach to remove material to make the cavity. Two adjacent negative features may form a stacked-up depth. Similarly, placing a negative feature inside another negative one will also form a stacked-up depth. Mold shape characteristics are identified from the part shape by feature mapping and model refinement discussed in the next section.

SIGNIFICANT ITEM EXTRACTION AND MODEL REFINEMENT

The identification of global geometric characteristics and the extraction of significant items of the characteristics is performed as a refinement procedure (see Figure 2) where high-level feature relationships are first identified to support the specialization for the geometric characteristics of individual features as well as the detailed spatial relationships among features. The significant items formed from feature interactions are then computed or derived based on the detailed spatial relationships and geometric details.

Based on a feature-based part model, features located on the same plane are identified first. Spatial relationships - "is in", "adjacent to" and "coplanar" between these features and their geometric details are then derived. The "is in" and "adjacent to" relationships are identified by searching edges shared by two features in the boundary representation model of a part, and distinguished based on the type (i.e. positive or negative) of feature involved.

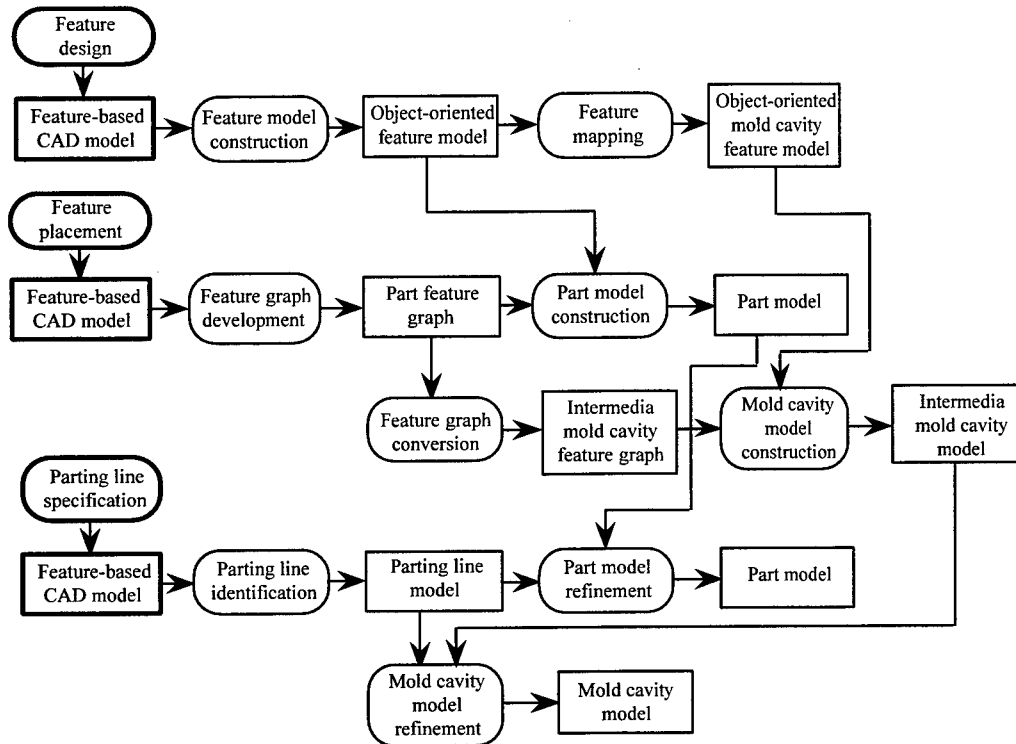


Fig. 2. Procedure for model refinement.

Identification of shape characteristics is seen as "specialization", i.e., to specialize high-level semantics - "is in", "adjacent to" and "coplanar" - into specific cases for particular application purposes. Three types of cavities - "hole", "groove" and "step" can be extracted from the "is in" relationship. Derivation of cavity type is based on the number of feasible directions for placing a negative feature to create a cavity. According to the cavity type, significant items, such as "depth" and "cross sectional area" of a "hole", "depth" and "width" of a "groove", and "depth" of a "step", can be extracted from geometric details.

The purpose of feature mapping is to convert a part feature into a mold cavity feature to support mold-related evaluations. A protrusion feature will become two depression features in the main cavity, if it is split by parting lines, otherwise, it maps into a depression feature. A depression feature can map into one protrusion feature and two depression features if it is not split by parting surfaces. On the other hand, the mapping results depend upon the location of the parting surface if it is split by parting surfaces. Since a parting surface perpendicular to the section plane of the feature will cause an undercut, an insert may be required. The mapping result in this case is two depression features and an insert. On the contrary, the mapping result will be two depression features if the parting surface is parallel to the section plane of the feature. A secondary feature will map into one or two complementary secondary feature(s) depending on whether it is split by the parting surface or not.

SYSTEM IMPLEMENTATION

The computer-aided integrated design for injection molding is implemented using Visual C++ in a WindowsNT™ environment on an Acer ALTOS 9000™ PC-server networked with 6 PC-clients, located in the Computer-Aided Concurrent Engineering Research Lab at National Cheng Kung University, Taiwan, ROC. A commercial CAD system (Pro/ENGINEER™ from Parametric Corporation), a commercial expert

system shell (Nexpert/ObjectTM from Neuron Data Corporation) and a commercial database system (AccessTM from Micro Soft Corporation) were selected as software components for system development.

There are three environments in the system: a design environment, a knowledge-based environment and a database environment. The design environment, consisting of a feature design module, a preliminary design module, a parting line/surface design module and a detail design module, is supported by a knowledge based environment, a geometric reasoner and the database environment. This system is fully integrated with a mold design system, a mold manufacturing process planning system and an injection molding process design system, which forms an environment for concurrent molding product and process development [9]. The knowledge-based environment contains a product design knowledge base, a mold design knowledge base and a process planning knowledge base. Besides supporting product design, the mold design and process planning knowledge bases also support a mold design environment and a computer-aided mold manufacturing process-planning environment respectively.

CONCLUSION

This paper presented an integrated design for injection molding framework that accommodates functionality, producibility and cost effective concerns. Based on this framework, a computer-aided integrated design for injection molding system was developed to guide the molding product design process and aid developers to improve development speed, consistency, and accuracy.

The major practical results of this work include:

1. A framework and method for integrated product design for injection molding. This gives a better way to design a product to fulfill the requirements of functionality, cost effectiveness and manufacturability.
2. A systematic approach to the development of a computer-aided integrated design for an injection molding system. This approach can be used to help develop of other computer-based tools or systems.
3. A computer-based integrated design for an injection molding system that provides consistency and systematic analysis and design of the shape of molding components. The results of this research will help to rationalize and automate molding product designs and improve the efficiency, quality and reduce the cost of product development.

ACKNOWLEDGEMENT

This research is funded by National Science Council, Taiwan, ROC under Grants: NSC87-2212-E-006-002, NSC87-2212-E-006-001, and NSC-87-2212-E-151-001.

REFERENCES

1. D.V. Rosato, P.E. Rosato, 1986. Injection Molding Handbook. Van Nostrand Reinhold Company, NY.
2. C. J. Haddad, 1996. Operationalizing the Concept of Concurrent Engineering: A Case Study from the U.S. Auto Industry. IEEE Transactions On Engineering Management, 43(2), 124-132.
3. B. Prasad, 1996. Concurrent Engineering Fundamentals: Integrated Product and Process Organization. Prentice Hall, New Jersey.
4. R.P. Smith, 1997. The Historical Roots of Concurrent Engineering Fundamentals. IEEE Transactions On Engineering Management. 44(1), 67-78.
5. G. Boothroyd, P. Dewhurst, W. Knight, 1994. Product Design for Manufacture and Assembly. Marcel Dekker, Inc. New York, USA.
6. D.C. Anderson, T.C. Chang, 1990. Geometric Reasoning in Feature-Based Design and Process Planning. Computer and Graphics, 14(2), 225-235.
7. J. R. Rossignac, 1990. Issues On Feature-Based Editing and Interrogation of Solid Models. Computer & Graphics, 14(2), 149-172.
8. Y.-M. Chen, C.-L. Wei, 1997, Computer-Aided Feature-Based Design for Net Shape Manufacturing. Computer Integrated Manufacturing Systems, 10(2), 147-164.
9. Y.-M. Chen, 1997, Development of a Computer-Aided Concurrent Net Shape Product and Process Development Environment. Robotics and Computer Integrated Manufacturing, 13(4), 337-360.

Artificial Psychology – An Attainable Scientific Research on the Human Brain

Zhiliang Wang, Lun Xie

School of Information, University of Science and Technology (USTB),
Beijing, 100083, P.R.China
Email: wzl@public.bta.net.cn

ABSTRACT

This paper presents a novel original theory of artificial psychology based on artificial intelligence. We analyze human psychology comprehensively from the context of information science research methods, especially in regard to aspects of emotion, willingness, character, creativity and the realization of artificial machines. We also present tentative definition, purpose, rules, research content, application domain and algorithms, in order to establish an architecture for artificial psychology and promote artificial intelligence research to probe deeper into the human mind and so, to generate higher levels of development in the field.

Keywords: Artificial Psychology, Cybernetics, Software

INTRODUCTION

(Theoretical Foundation, International/Domestic Trends, Applications, Perspectives, R&D Values)

At this most important period in time, (the turn of the 21st century), many senior researchers are delving into the problem of how to combine life sciences with information science to help promote both areas to their full potential. Most research focuses on electronic information imitation of the human brain, specific aspects of human intelligence and human behavior. This paper presents some probing methods and tries to develop a novel research field in automation science or, especially, in information science, on the basis of human brain science, psychology, neural science, and automation theory.

It is well known that on the one hand, research into Artificial Intelligence has been done to a very high level. But on the other hand, its purpose has been limited to imitating human intelligence talents such as judging, inferencing, proving, explaining, identifying, sensing, planning, learning and problem-solving. The main tasks of these activities is how to present, acquire and utilize knowledge. This is not nearly a wide-enough research range since human psychology and related activities include not only sensing, believing, memorizing and thinking, but also emotion, feelings, willingness, character and creativity. The core of this paper is aimed at a comprehensive understanding of human psychology, (especially feelings, willingness, character, and creativity) by means of a systematic integration of psychology, brain science, information engineering, computation science and novel theories of automation science.

The most valuable research work in Artificial Psychology are the comprehensive post-imitation and machine realization (computation and algorithmic) being done on human psychological activities. Application perspectives include development of robots with colorful feelings and intelligence, virtual mechanics, and control systems based on the human brain (feelings + sensing = behavior). There is a vast potential to apply artificial psychology theory to the design of more humanized commodities and to help set up man-to-man and man-to-machine social environments.

Although, no one has previously devised a theory of Artificial Psychology, there is a trend towards R&D in AP, especially in Japan. According to the explanation in Japanese academia, KANSEI Engineering is the integration of feelings, emotion, sense, and energy. In order to maintain a leading role in information technology into the 21st century, various experts under the organization of the Japanese Economic Ministry, have been researching the project of "human media" since 1995 with support of both the Economic and Education Ministries in Japan. The results presented in the Proceedings of the International Conference held in Nagoya in 1997 and some electronic academic materials are just preliminary. (see References).

For example, the definition of KANSEI Engineering is extremely unclear and its application field is just one aspect of human psychological activities (sensing), so, this paper presents a more general theory to study how we can imitate human psychological activities more comprehensively. The American Government has taken research into this Science most seriously in recent years and has passed declarations in both the US Senate and the US House of Representatives to define the code between 1990~2000 as the Decade of the Brain which epitomizes the booming activities taking place in Brain Science. Human sensing, movement learning, memory, thinking, emotion and behavior are all functions of the Brain, which is to say, that human psychology is the function of the Brain.

Brain Science is a biological science, which belongs to natural sciences. Psychology is the image of real facts in the brain, which is, in reality, a social science. Brain Science is the foundation of psychology and psychologists must be familiar with the structure and function of the human brain. So, brain science and psychology are the main theoretical bases for artificial psychology. Since imitation of human emotion by means of a machine is not recognized in the field of Artificial Intelligence in China, there is hardly any research work being done in this area. However, we deem it advisable that the imitation of human emotion by means of computers is one of the forward-technologies in the field of information science, which is controlled by the National High-Tech Planning Structure as a feature group within the National Science Committee of China. So, it is necessary and urgent for us to develop and keep abreast of the latest international developments in this newly-emerging field. Above all, on the basis of an extensive analysis of recent trends in the relevant academic areas, this paper presents the essential forward-scientific characteristics of Artificial Psychology which includes not only theoretical values of great importance, but expansive application ranges as well.

ARTIFICIAL PSYCHOLOGY

It is well known that psychology is a basic science of psychological phenomenon generation, development and activities. Psychology is the integration of thought, emotion, memory, willingness, character, and creativity, i.e., it is image processing of real facts contained within the human brain. Artificial Intelligence is a kind of science based on knowledge whose academic domain is theoretical knowledge (advanced progress of psychological activities), while artificial psychology deals with lower processes of psychological additives whose academic domain is the psychology of the logical activities of emotion, willingness, character, creativity with much more fuzzy performance. The relationship between artificial psychology and artificial intelligence can be seen in Fig. 1. In our opinion, artificial psychology, which is based on artificial intelligence, is an advanced stage of artificial intelligence with more extensive content, both are inter-related and complemented.

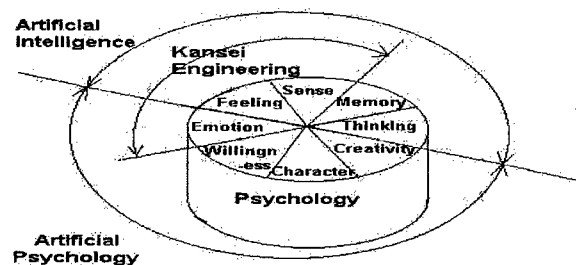


Fig. 1. Relationship Between A.I. and A.P.

From Fig. 1, we can see that the theory of Artificial Psychology is different from KANSEI Engineering as presented by Japanese academics. Engineering Psychology in the Psychological field and the harmonic interface of man-machine as promulgated by the National Natural Science Committee in China are also a subset since the application of Artificial Psychology must be more expansive and novel. The paper tries to setup the theoretical architecture of Artificial Psychology (rules, purposes, applications, methods) and to realize its practical application to develop a novel academic field in automation and/or information science. Artificial Psychology is an interdisciplinary science as shown in Fig.2. Its foundation stems from Brain Science, Engineering, Kansei Engineering, Linguistics, Law, Information Science, Automation and Artificial Intelligence. Its practical application includes technical support of emotional robots, humanizing

commodity design, sensible market development, Artificial Psychology Programming Languages, Artificial Creative Techniques, Virtual Techniques for Human Psychology Databases and Mathematical Models, Harmonic Interface for man-machine and IS multi-channel Interfaces.

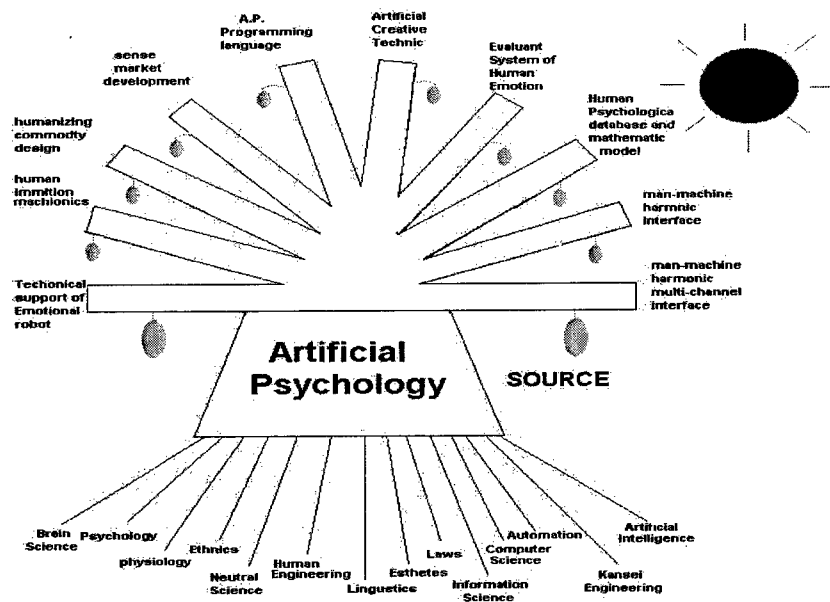


Fig. 2. Artificial Psychology- an Interdisciplinary Science.

The theoretical structure of artificial psychology is given in the following summary:

1. **Definition:** The realization of an artificial machine that contains computer model algorithms that comprehensively mimic human psychological activities including emotion, willingness, character, creativity mainly by means of information techniques.
2. **Purpose:** To build up a happy harmonic social environment of man-to-man, man-to-machine and man-to-nature interactions.
3. **Rules:**
 - 1) Imitate human psychology actively.
 - 2) Virtue, esthetic happiness are the essential rules.
 - 3) Creative imitation of human psychology
 - 4) Machine serves as a means to mimic man forever.
4. **Aims:**
 - 1) To present the notion of artificial psychology by means of the theoretical foundations of A.I. with the integration of psychology, neuroscience, neural science, information science, computation, and automation;
 - 2) To imitate human psychological activities (emotion, willingness, character, creativity, etc.) comprehensively, by setting up the theoretical architecture of Artificial Psychology (purpose, rules, applications, methods); and
 - 3) To fulfill the practical application of A.P.
5. **Research needs and vital problem solutions:**
 - 1) To set up the theoretical architecture of Artificial Psychology, especially to overcome the limitation of definitions, rules, and content, in order to apply the moral rules of human beings which is not currently included in Artificial Intelligence.
 - 2) To set up the theoretical architecture of Artificial Psychology by means of the theoretical achievements of A.I. and the relationships between A.I. and A.P..
 - 3) To develop machine algorithms to repress bad emotions, which is decided by the rules of A.P.
 - 4) To mathematically scale human psychological information whose main achievements were made by Japan academics.
 - 5) To set up a control mode of sense, feeling and emotion as the decision of behavior to imitate

the control mode of human brain. The process goes from sensing (instruments) to feeling (analysing) to action (behavior) with emotion (adaptation modules) over-riding the normal feedback or feedforward control. (See Fig. 3.)

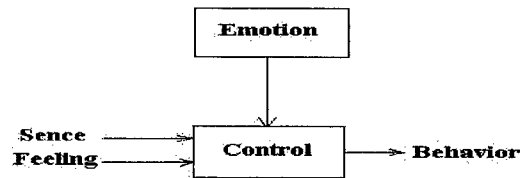


Fig. 3. The Emotional Control of the Human Brain.

- 6) To probe the building-up of a programming language for Artificial Psychology which is a challenging task. The programming language of A.I. is the presentation of knowledge and logical inference. In A.P., the programming language must be a kind of associative language whose character is associative inference, chaotic computation, divergent thinking and fuzzy induction. The comparison between these is shown in Fig.4.

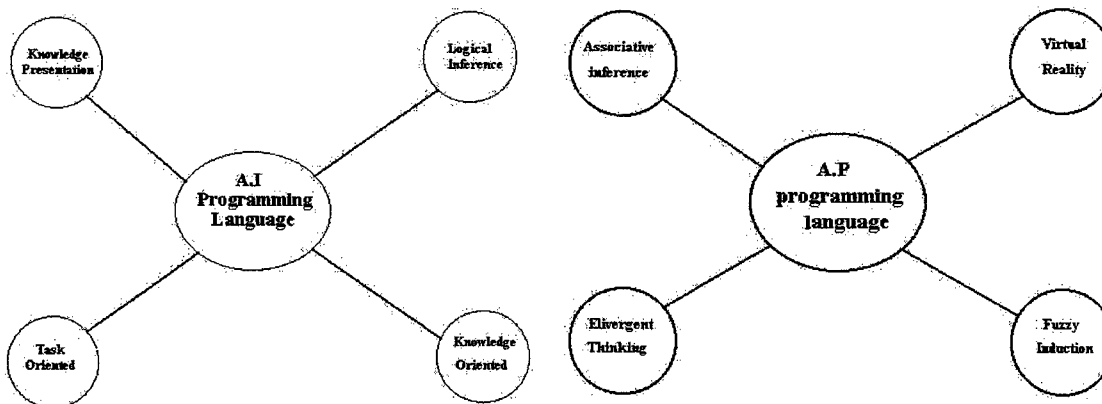


Fig. 4. A.I. and A.P. Programming Languages.

- 7) Computation algorithm for emotion cultivation.
 - 8) **QiGong** - the building up of psychological implication and functions of a human being.
 - 9) Machine realization of inspiration scintillation as depicted in Fig. 4.
- 6) **Research Methods & Technologies**
- 1) **Methods:**
Set up the theory of AP by means of fuzzy logic, neural networks, chaotic theory, advanced algorithms, on the basic of psychology and brain science.
 - 2) **Technologies:**
By researching psychology, brain science, with reference to Kansei Engineering in Japan, we will adopt the theory of fuzzy-logic, neural networks and chaos theory to induce a mathematical model of human psychology and realize the theory of A.P. into software packages such as A.P. databases and process modules which "sense" information and to evaluate and imitate human psychological status by means of virtual reality to create a psychological evaluating "brain" to promote development of A.P.
- 7) **The Prospects of Application and Achievements**
It is believed that once a theoretical architecture of A.P. has been set up and human psychological information has been scaled, A.P. can be realized on a computer such as imitation of inspiration scintillation and the application of emotions in functional decision-making. The main academic areas focus on application of emotional control theory in robots; application of sense information processing in commodity design and market development and application of inspiration in creative engineering.

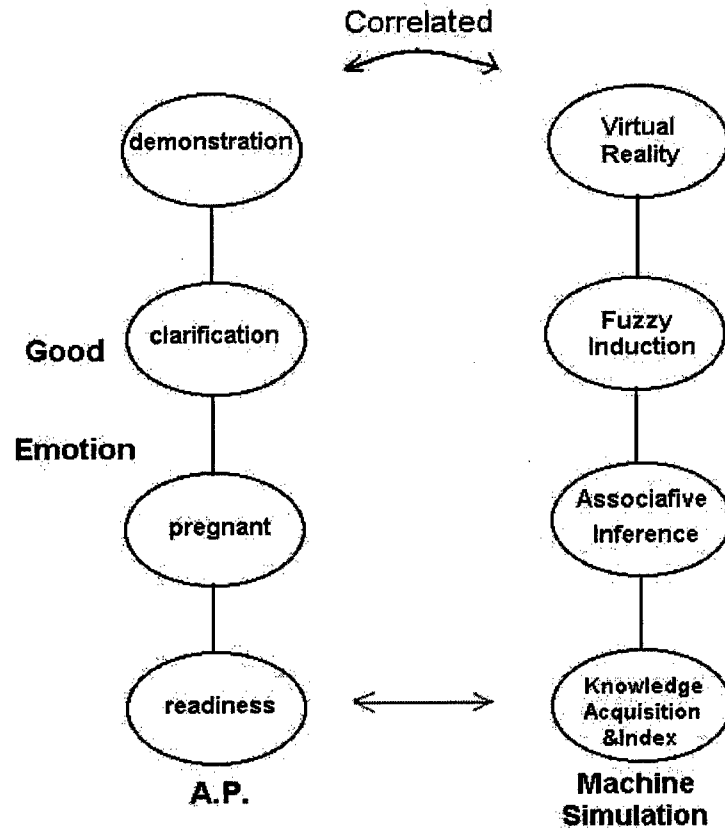


Fig. 4. Machine Realization Scintillation.

CONCLUSION

The theory of A.P. as presented in this paper is a novel notion, still in its infancy. It is true that there are many complicated problems to be resolved in A.P. such as belief factors of artificial imitation of natural psychology; fuzziness; chaotic psychological activities; accuracy of psychological models; and the importance of adaptability of A.P. to human morals. By an interdisciplinary academic collaboration of the fields of information science and life science presented in this paper, the perspective and applications are extensive and far-reaching. This paper is a preface to the theory of A.P. We invite experts and science researchers all over the world to join us to develop potential myths and markets in A.P. in the future.

REFERENCES

1. Qian Ye, D.A. Fu., 1996. Fuzzy reasoning model on feeling and action, Japan Human Eng., 5, 230-238.
2. Jia Teng Junyi, 1998. Kansei Agent and Human Media Database -- Kansei working place. system, controlling and information, 5, 253-259.
3. Chang Gu Chuanlong, 1997. The fundamental conception of multi-media Kansei composing processing and the assessment of its testing system, Thesis Japan Information Processing Institute, 8, 1517-1530.
4. _____, 1998. Speciality: Kansei and language information processing -- expressing Kansei and language through computers, Japan Computer Today, 1, 4-45.
5. Da Tsiao Li, 1996. Specific integration: fusion of science and arts - expectations of Kansei Eng., I.IEE, 1, 4-7.
6. X.S. Youxiang, 1997. Evaluation of handkerchief design based on Kansei Eng., I.IEE, 117-c7, 934-939.
7. A. Buzhen, 1997. Composing music from landscape picture with fuzzy method, I.IEE 117-c4, 452-457.
8. Zou Zhi, 1985. Mathematical arts. Attempt I: cheng cheng literature, 111, 106-126.
9. Zou Zhi, 1986. Mathematical arts. Attempt II, cheng cheng literature, 117, 80-101.
10. Tu Zyo Shangzi, 1997. Art and technology, J. Fuzzy Inst. Japan: Considering Kansei Eng., 5, 648-656.
11. Zuo Zhi Tsingfu, 1997. Science and Kansei. J. Fuzzy Inst. Japan: Considering Kansei Eng., 6, 861-869.
12. Ba Mu Shaohong, 1997. Kansei physical measurement, J. Fuzzy Inst. Japan, 3, 318-326.

13. Jiu Jin JianYong, 1998. Word processor with Kansei Retrieval. *J Fuzzy Inst. Japan*, 5, 318-326
14. TsingShui Yixiong, 1998. Informationalized society and Kansei Eng., *J. Fuzzy Inst. Japan*, 5, 804-812.
15. Matsuda, N.; Nakamura, K., 1997. Interactive support for decision making, in *Design of Computing Systems: Cognitive Considerations*. Proc. 7th Inter. Conf. on HCI., San Francisco, 24, 479-82.
16. Hayashi, T.; Hagiwara, M., 1998. Image query by impression words - the IQI system, *IEEE Trans. on Consumer Electronics*, 44(2), 347-52.
17. Yamaoka, T., 1997. A new design concept and method based on ergonomics and Kansei engineering, in *Design of Computing Systems: Cognitive Considerations*, Proc. 7th Inter. Conf. on Human-Computer Interaction, San Francisco, 2, 547-550.
18. Hata, S.; Miyashita, Y.; Hanafusa, H., 1997. Color sensitivity of human in visual inspection, *IEEE/ASME Inter. Conf. on Advanced Intelligent Mechatronics '97*.
19. Shibuya, K.; Chikaoka, Y.; Koyama, T.; Sugano, S., 1997. The planning of violin playing robot with KANSEI information-algorithm to decide bowing parameters from timbre. Proc. 6th IEEE Inter. Workshop on Robot and Human Communication. RO-MAN '97, Sendai, 230-235.
20. Shibuya, K.; Sugano, S., 1997. Human motion planning in violin playing using KANSEI, 1997 IEEE SMC Conf., Orlando, 3, 2638-2643.
21. Hayashi, T.; Hagiwara, M., 1997. An image retrieval system to estimate impression words from images using a neural network. 1997 IEEE SMC Conf., Orlando, 1, 150-155.
22. Tanaka, S.; Inoue, M.; Ishiwaka, M.; Inoue, S., 1997. A method for extracting and analyzing 'Kansei' factors from pictures, 1997 IEEE 1st Workshop on Multimedia Signal Processing, Princeton, 251-256.
23. Inoue, S.; Ishiwaka, M.; Tanaka, S.; Park, J.-I., 1997. An Image Expression Room, Proc. Inter. Conf. on Virtual Systems and MultiMedia, VSMM '97, Geneva, 178-87.
24. Bianchi, N.; Berthouze, L., 1997. Supervised self-organization of user's Kansei model for image retrieving, Proc. 2nd Inter. Conf. on Cogn. Tech., (No.97TB100146), Aizu-Wakamatsu City 185-189.
25. Koshimizu, H., 1997. Future trends of visual inspection systems. a prospect for neo-machine vision applications, QCAV 97, Inter. Conf. Quality Control by Artificial Vision, Le Creusot, France, 36-40.
26. Tamano, K., 1996. Optical method for processing information depicted by picture, Proc. 4th Inter. Conf. on Soft Computing, Fukuoka, Japan, 2, 622-625.
27. Shibata, T., 1996. Artificial emotional creature project for intelligent system-human robot interaction, Proc. 4th Inter. Conference on Soft Computing, Fukuoka, Japan, 1, 43-48.
28. Tsuchiya, T.; Matsubara, Y.; Nagamachi, M., 1995. A study of Kansei rule generation using genetic algorithm, Proc. 6th Inter. Conf. on Human-Computer Interactions, Tokyo, 191-196.
29. Matsubara, Y.; Nagamachi, M., 1995. Hybrid Kansei engineering system and design support, Proc. 6th Inter. Conf. on Human-Computer Interactions, Tokyo, 161-166.
30. Shibata, Y.; Fukuda, M.; Katsumoto, M. A hypermedia-based design image database system using a perceptual link method. *J. Management Information Systems*, 13.
31. Kato, T.; Hirai, S.; Tomita, F.; Niki, K.; Higuchi, T., 1996. Human media technology for future information environment-kansei media technology and its application to human communication. *Bulletin of the Electrotechnical Laboratory*, 60(8), 13-49. (Japanese).
32. Yamaguchi, M.K.; Kato, T.; Akamatsu, S. Relationships between physical traits and subjective impressions of faces - age and gender information. *Electro. and Comm. in Japan*, Pt 3, 79(10), 23-34.
33. Kondo, T.; Yamagiwa, T.; Yamanaka, K.; Yamamoto, M., 1996. Detection of skill in skiing by motion analysis, *Trans. Inst. Electronics, Information and Communication Eng., Scripta Technica*.
34. Nomura, J., 1994. Virtual reality technologies and its applications to industrial use, *Virtual Reality Software and Technology*, Proc. VRST '94 Conf., Singapore, 125-42.
35. Kwon, K.S.; Lee, S.Y., 1993. A study on the product design considering human emotion as a part of human interface technology, Proc. 32nd SICE Conf., (IEEE 93 TH 0575-1), Kanazawa, Jap., 1091-1093.
36. Ishihara, S.; Ishihara, K.; Matsubara, Y.; Nagamachi, M., 1994. Self-organizing neural networks in Kansei engineering expert system, ECAI 94, 11th European Conf. on A.I., Amsterdam, 231-5.
37. Kawakami, F.; Morishima, S.; Yamada, F.; Harashima, R., 1994. Construction of 3-D emotion space based on parameterized faces, Proc. 3rd IEEE Inter. Workshop on Robot and Human Communication. RO-MAN '94, Nagoya (Cat. No.94TH0679-1) 216-21.
38. Enomoto, N.; Nagamachi, M.; Nomura, J.; Sawada, K., 1993. Virtual Kitchen System using Kansei Engineering, Proc. 5th Inter. Conf. on Human-Computer Interaction, Orlando, 2, 657-662.

Soft-Object Technology for Flexible Machining Systems (FMS)

Ahmed Hambaba

CISE Department, College of Engineering
San Jose State University, San Jose, CA 95192-0180
Tel: (408) 924 3959 Email: Ahmed_Hambaba/SJSU@sjsu.edu

ABSTRACT

Individual manufacturing sites are expected to become more demand-driven than plan-driven. Manufacturing industries are faced with the challenge of responding more rapidly to changing markets and evolving business opportunities. Focus of factory-floor information systems is on the operation of production equipment and the control of processes. There is no direct or regular communication with design and engineering systems. As a result, upstream information systems are unaware of factory-floor information details, such as the status of work-in-progress (WIP); records of past process performance; availability of appropriate tools, labor, and materials. Middle-level information systems, known as manufacturing execution systems (MES), bridge this critical information gap between upstream and downstream activities.

INTRODUCTION

Manufacturing environments are moving away from stable, high volume environments to more flexible, lean environments. Computer science has moved away from the central mainframe type model to a more modular client/server-distributed object model. These changes were recognized by Texas Instruments (TI) and the Department of Defense, who began a joint project in Oct. 1988 known as the Microelectronics Manufacturing Science and Technology (MMST) program. [3]

The purpose of the MMST was to revolutionize semiconductor manufacturing. Their objectives included improvement of the manufacturing process and the creation of a next generation CIM system architecture that could exploit new manufacturing technologies while addressing shortcomings in current CIM technology [3]. The MMST implemented a CIM architecture using Object-Oriented technology, which resulted in reduced development time, improved interfaces, improved module reusability, improved flexibility and improved extendibility. The software was so reusable that the embedded machine control in a TI machine was moved to another supplier's machine with no changes to the reused classes whatsoever [3].

The CIM implementation resulted in almost complete elimination of operational errors due to the extended role that the system played in running the factory [3]. This CIM system, along with the improvements made in the manufacturing process, reduced operational cycle time from the typical 30 days to 3 days for single wafer lots.

SEMATECH, a consortium of semiconductor manufacturing companies dedicated to creation and implementation of standards, followed the lead established by TI and began creating an object oriented CIM framework specification based on the MMST. SEMATECH calls it the Computer Integrated Manufacturing (CIM) Framework Specification.

The SEMATECH CIM specification describes the framework of a computer system that will "assist in the creation of an integrated, common, flexible, modular object model leading to an open, multi-supplier CIM system environment". [1]

The SEMATECH CIM specification limits itself to the description of manufacturing software on the manufacturing execution systems (MES) level. MES is that category of software that is between direct equipment control and corporate planning [2], as shown below in Figure 1.

The SEMATECH framework is a set of functional components designed to work together to form an integrated manufacturing system. The components covered in this specification are listed in Table 1.

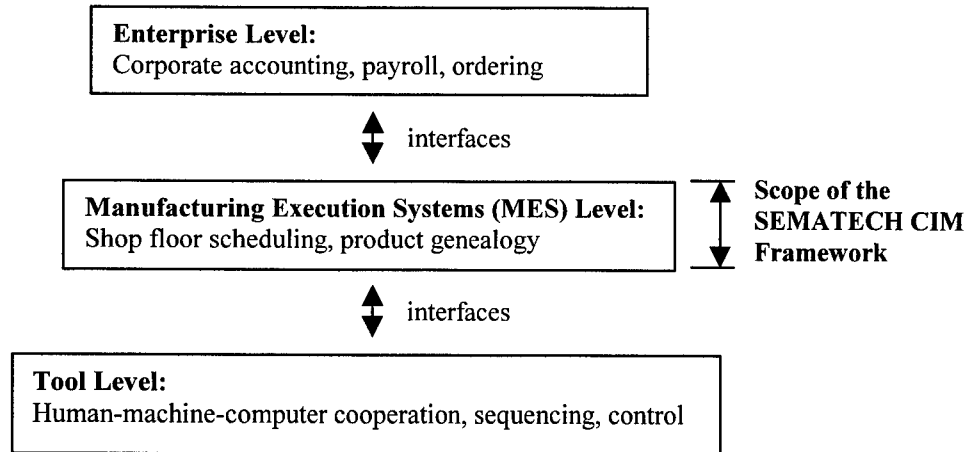


Fig. 1. Enterprise System Context [2].

Table 1: SEMATECH CIM Framework Specification Components [1].

Factory Services	Material Management
<ul style="list-style-type: none"> Document Management Version Management History Management Event Broker 	<ul style="list-style-type: none"> Product Management Durable Management Consumable Management Inventory Region Product Specification Bill of Material
Factory Management	Material Movement
<ul style="list-style-type: none"> Factory Product Release Factory Operations Product Request 	<ul style="list-style-type: none"> Material Movement
Factory Labor	Advance Process Control
<ul style="list-style-type: none"> Person Management Skill Management 	<ul style="list-style-type: none"> Plug In Management Plug In Execution Control Management Control Execution Control Database Data Collection Plan
Machine Control	Process Specification Management
<ul style="list-style-type: none"> Machine Management Recipe Management Resource Tracking 	<ul style="list-style-type: none"> Process Specification Process Capability
Schedule Management	
<ul style="list-style-type: none"> Dispatching 	

CIM SPECIFICATION

This section will review the components of the specification that pertain the most to machine control and work management. These components will be reviewed as they were written in the specification. Later sections will review the actual implementation strategy, as it does differ from the specification. The components that make up the structure will be discussed first, followed by the Job flow structure.

Components

The two main component groups we will be reviewing are the Machine Group and the Factory Component. Their relationship is illustrated in the following diagram, Figure 2.

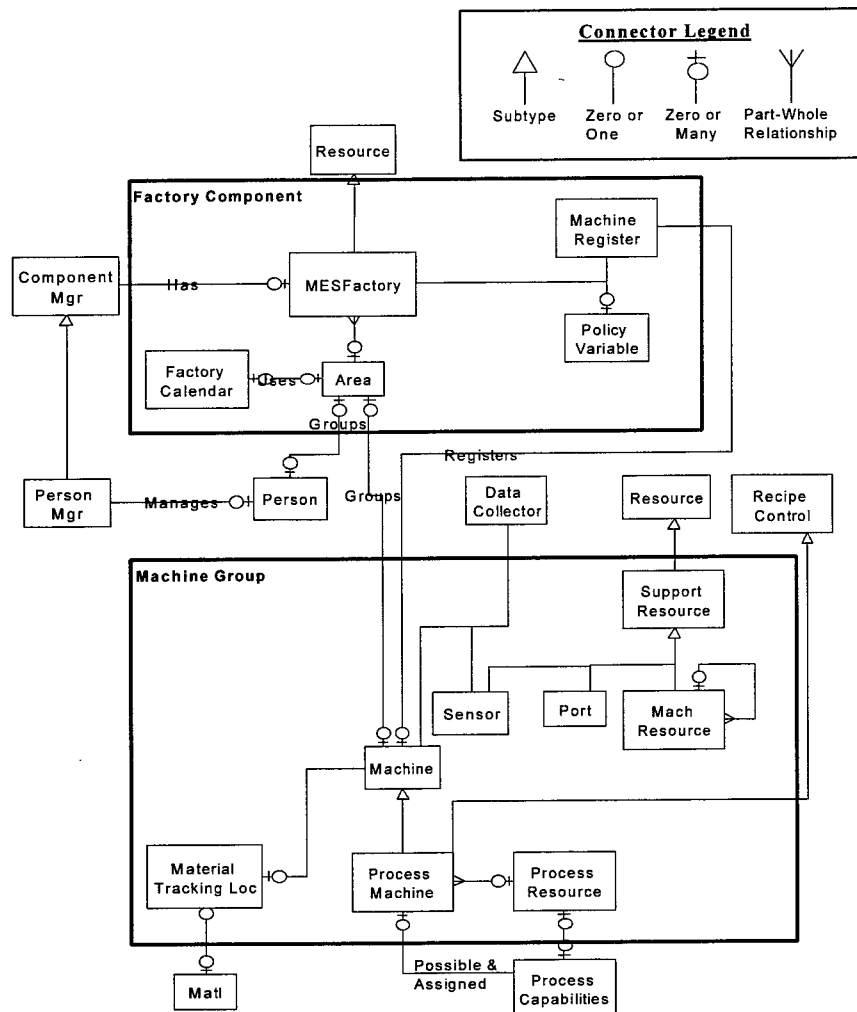


Fig. 2. SEMATECH CIM Specification Components Diagram [1].

The Machine Group is made up of ProcessMachine, Sensor, Port, ProcessResource, and none or many Material Tracking Locations:

- The ProcessMachine represents a device for processing product.
- The ProcessResource is a subinterface of the ProcessMachine. There can be more than one ProcessResource per ProcessMachine as in the case of cluster tools. The ProcessResource performs the Process Operation.
- The Port represents an access point to the Machine for material.
- Material Tracking Locations are places where material may be held.
- Sensor represents a data collection point for the Machine.[1]

The Factory Component is made up of a MESFactory, one or more Areas, a Machine Register with none or many Policy Variables, and a Factory Calendar. The objects are described below:

- The MESFactory represents the main interface to the Factory Component. It is a composite object that represents all factory resources.
- Areas represent groups of resources, grouped either logically or geographically. This is an optional component.
- The Factory Calendar is used to keep track of all Factory scheduling, such as holidays, shut-downs, etc. This represents each day in the life of a Factory.

- The Policy Variable is used to add additional configurability to the factory, by providing a vehicle for more advanced business rules.
- The Machine Register interface maintains a list of all machines known to the factory. There is only one Machine Register per factory.

Figure 2 also illustrates how the Machine Group relates to the Factory Component through the Area and the Machine Register. Each Factory can have zero or more Areas, which can be either functional or geographical groupings. Machines are grouped in these Areas, and the Factory can query its areas for the Machines that reside in them. In order to determine which machines are available to the Factory to do work, the Factory uses the Machine Register. As Machines come up they register themselves with the Machine Register. The Factory queries the Machine Register to learn what Machines are available and what their current state is.

This illustrates the Factory structure, but not the flow of work. For that, we need to examine the Job flow structure.

Job Flow Structure

Work is passed into and through the factory in the form of Jobs. A Job is defined as a unit of work that is requested of and performed by the factory. The flow of work is defined by the following diagram.

Figure 3 illustrates how work requests make their way from the top level of the MES, the Enterprise, as a Product Request, down to the lowest level, the ProcessMachine, as a ProcessMachineJob. This process follows these steps:

1. The Enterprise works with the Product Request Manager to create a Product Request.
2. The Product Request is broken into Lots Jobs by the Product Request Manager, which then releases the Lots Jobs into the Factory in an order that maximizes the Factory's resources.
3. The Factory works with Dispatcher and the Scheduler to determine the best schedule of activities for the factory.
4. The LotJob Executor then uses this information to request work from the ProcessMachines.
5. The ProcessMachines control floor level machine actuation to get the Product created.
6. The ProcessMachine informs the Factory when the Job is complete, the Factory informs the Product

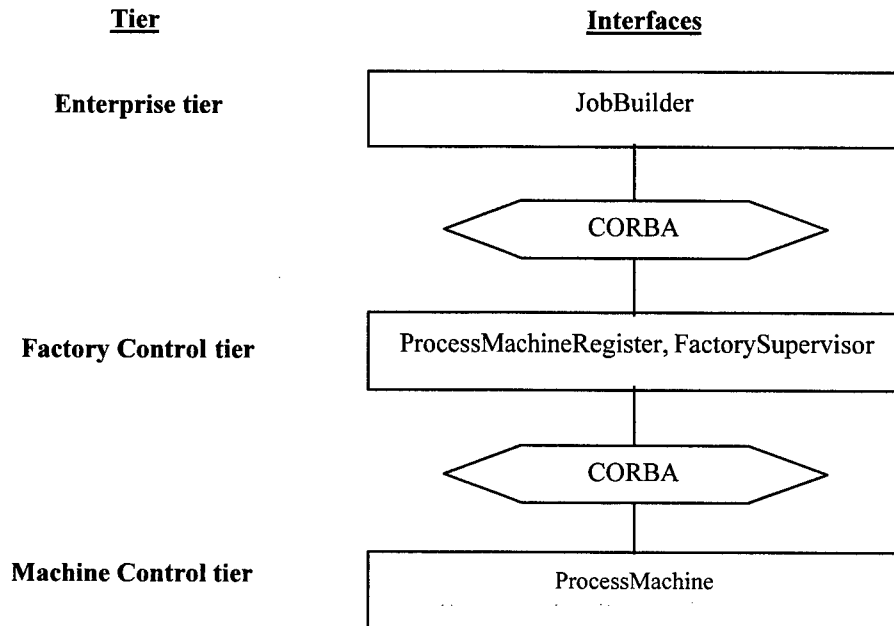


Fig. 3. Implementation Tiers.

When the LotJob is complete, Request Manager and the Product Request Manager inform the Enterprise when the Product Request is complete.

The next section will review the design of this project implementation and how it differs from the specification described above.

SYSTEM DESIGN

Architecture

This project implementation has three main tiers: the Enterprise tier; the Factory tier; and the Machine Control tier. These tiers are separated by the CORBA ORB, described in the Technology Used section, and are represented by the interfaces shown in Figure 3 below.

The Machine Control tier is the bottom tier, which consists of one or more ProcessMachines and represents and actuates the physical machines in the system. The Factory tier is the middle level that is made up of the FactorySupervisor and ProcessMachineRegister. The Factory tier accepts jobs on behalf of the Factory and delegates subjobs to the Machine Control tier. On the top is the Enterprise tier, which uses the interface EnterpriseExecutor to request jobs from the Factory tier and to represent the interface that an Enterprise system would have with the factory.

Technology Used

The technologies used were the Java Developers Kit (JDK) 1.1.7B, Visigenic Visibroker CORBA Compliant ORB v.3.1, Symantec Visual Café Database Development Edition v.2.5a, and WinEdit v96W.

JDK 1.1.7B is the Java compiler and virtual machine created by Sun Microsystems. It is assumed that the reader is familiar with the Java language.

Symantec Visual Café Database Development Edition v.2.5a is a Java Integrated Development Environment (IDE) that allows for drag-and-drop creation of GUI layouts. Visual Café has many other features like a debugger but was used in this implementation solely for GUI development.

WinEdit v96W is an enhanced text editor program. This editor offers such features as different color fonts for different parts of code (this greatly increases readability) and a compiler and debugger. For this project WinEdit was used only for editing source code. All compiling, running, and debugging was done on the JDK command line.

Visigenic Visibroker v3.1 is a CORBA compliant ORB (CORBA is explained in the next section).

CORBA

The Object Management Group (OMG), a consortium of companies from all facets of the computer industry, has defined the Common Objects Request Broker Architecture (CORBA). CORBA is the definition of middleware that allows intelligent components to discover each other and interoperate on an object bus [4]. This object bus is referred to as the Object Request Broker (ORB). The ORB abstracts away the information needed for remote components to communicate. Components can interact with each other across networks, servers, and operating systems as if they all were residing on the same machine. This makes it very easy for the developer to create networked client/server applications.

The component's boundaries are defined using Interface Definition Language (IDL). The IDL file is then compiled to generate client side stubs and server side skeletons, which permit components to be accessible across languages, tools, operating systems, and networks. This allows the developer to worry about object interaction rather than all the activities required to locate remote objects and to pass information across a wire. CORBA also defines an extensive set of bus-related services for creating and deleting objects, accessing them by name, storing them in persistent stores, externalizing their states, and defining ad hoc relationships between them. [4] The main CORBA service used in this implementation was the Naming service.

The Interfaces

The project implementation diagram reveals the need to study both abstract interfaces and non-abstract interfaces. The abstract interfaces include Resource, JobRequestor, JobSupervisor, and Job. The non-abstract interfaces are ProcessMachine, ProcessMachineJob, ProcessMachineRegister, FactorySupervisor, LotJob, and EnterpriseExecutor. The abstract interfaces are examined first.

Prior to discussing the interfaces, it is worthwhile to discuss the concept of Capabilities, or JobCapabilities as they sometimes appear in the interfaces. Capabilities are a fundamental part of this implementation and can be defined as a process step to be performed by a ProcessMachine. ProcessMachines maintain the Capabilities they can perform, and Jobs are made up of Capabilities to be performed. A JobRequestor will build lists of Capabilities and use them to request Jobs from JobSupervisors using these lists. The JobSupervisors will then verify that they can perform the Capabilities, and if they can, they return a Job to the JobRequestor. Capabilities show up throughout the following interfaces in the ways just described and other ways as well.

CONCLUSION

The CIM specification is a very complicated document. In attempts to isolate a certain set of functionality to implement, we came to the conclusion that this specification was written not to allow for free standing components, but rather to allow for components from different suppliers to be plugged into an infrastructure. The interfaces are very complex, and inheritance is many levels deep, leading to a situation where isolating pockets of functionality is almost impossible.

REFERENCES

1. Doscher, David, Editor, 1998. "Computer Integrated Manufacturing (CIM) Framework Specification Version 2.0", SEMATECH Technology Transfer 93061697J-ENG.
2. Freed, Ken, 1995. "Implementation Handbook for the Computer Integrated Manufacturing (CIM) Application Framework Specification 1.2", SEMATECH Technology Transfer 95092971A-ENG.
3. McGehee, John, et al., 1994. "The MMST Computer-Integrated Manufacturing System Framework", IEEE Transactions on Semiconductor Manufacturing, 7(2).
4. Orfali, Bob, Harkey, Dan and Edwards, Jeri, 1997, Instant CORBA, Wiley Computer Publishing, NY, NY.

A Monitoring Framework for Software Project Development

Ho-Leung Tsoi* and Derek Cheung**

* Software Quality Institute, Griffith University, Australia

** Division of Computer Studies, City University of Hong Kong, Hong Kong

ABSTRACT

Software project development includes a number of activities that result in a delivered product (software). As software becomes more and more expensive to develop, monitoring is an important task for project development and has been recognized as a difficult task in practice. There are a lot of unpredictable factors existing in the software development cycle that have become contributing factors to this problem. This paper gives an overview of the present state of the art of the software development projects. Moreover, a monitoring framework is proposed to help project management to get better understanding and make all activities running on schedule.

INTRODUCTION

Effective management is the ultimate objective of all software project managers. Although there are many methodologies and techniques available to be used for software development, many projects still suffer from late completion time and budget overruns. In 1984, a survey was conducted by Jenkins to address this problem [1]. The developers of 72 information system projects in 23 major American corporations were interviewed. Jenkins reported that the average effort and schedule overruns were 36% and 22% respectively. Researchers at the University of Arizona carried out a similar study and 191 responses were obtained [2,6]. They reported the average cost overrun was 33% -- very close to that reported by Jenkins.

Experience and findings from numerous studies have shown that software project management is difficult because there are a lot of uncontrollable factors existing in the development cycle. For instance, users may change the requirements at any time of the development cycle. All these changes may be a serious influence on the effort and duration of software development. As a result, process monitoring is a vital activity, required to be performed in software project management. It enables errors to be found early enough so that corrective action can be taken immediately. This paper presents a framework called *Software Monitoring Framework (SMF)* focused mainly on the way to monitor software development. A case study is presented.

PRINCIPLES OF MONITORING IN SOFTWARE PROJECT DEVELOPMENT

There is no doubt that project managers have to confront with many different kinds of problems, such as technical, management, and personnel, in the software development cycle. They are responsible for ensuring that all development processes are controlled appropriately in order to meet the budget and schedule. Unfortunately, software projects regularly get out of hand in the development cycle because there are a lot of uncontrollable factors exist in the development cycle.

For instance, users may change the requirements at any time of the development cycle. All these changes may be a serious influence on the effort and duration of software development. As a lot of factors can influence the development progress, process monitoring is necessary but has been considered difficult to accomplish. This paper presents a framework for monitoring the software development. Two principles are adopted:

Change exists in software development

Development of a software application is a dynamic environment characterized by many changing factors. Change may relate to technologies, personnel, or requirements. Some of these are as follows:

- Change in software size and complexity
- Change in development paradigm
- Change in development manpower
- Change in outside-support and/or inside-support

Software effort measurement must be performed dynamically

Since change exists in the development cycle, the measurement must be performed dynamically in order to collect the real-time information. In other words, the measurement is used to continuously keep track of the project progress of all activities to cope with change.

In summary, change exists in the software development cycle. The project managers need to continuously monitor the progress and take corrective actions if necessary. This paper proposes a process monitoring framework which helps the project managers to apply these two principles in the software development.

SOFTWARE MONITORING FRAMEWORK

To achieve a successful software system, all development processes must be continuously monitored. Thus they can be adjusted during development in order to cope with the problems of the changing world. This paper describes a framework to address the change problem during software development. As a result, the goals (work content) at different development phases can be achieved.

Figure 1 provides an outline, in a simplified form, of the framework developed in this paper. The stages of the system can be classified into (1) the acquisition phase and (2) the operation phase. Acquisition phase includes all the activities ranging from research and effort planning. Operation phase includes all activities during the actual software development cycle (Figure 2). The major tasks in the acquisition phase involve estimate the development effort, establishing effort management policy, establishing historical effort database, and identifying the critical development factors. The major tasks in the operation phase include reviewing monitoring and evaluation metrics, development effort monitoring, and evaluation collected information. The historical database developed in the acquisition phase can be used in the operation phase.

In summary, the acquisition phase is a pre-requisite for monitor and determination. The operation phase consists of monitoring mechanism, which collect information to make decisions and ensuring timely detection. Besides, a feedback channel, from operation phase to acquisition phase, is set up which operate via status and progress reports to compare actual progress with the plans based on the estimates. This framework ensures the timely monitoring and minimizing technical/management risks.

ACQUISITION PHASE

The acquisition phase includes the activities of estimating, planning, scheduling, budgeting, etc. The major tasks in this phase are as follows:

Estimation

It is commonly agreed that estimating is a prerequisite for good management. At the early stage of software development, the project management needs to determine the work-content. The work-content of the application can be determined by one of mathematical techniques or personal judgement. The result of the estimation let project managers to get the basic picture of the software application project. In summary, the effort estimation should satisfy three major requirements.

- Overall development effort will be identified.
- Be compatible with the data requirements for progress reporting.
- All critical cost drivers for the development process has been considered.

Assessing Development Factors

As previously mentioned, change will occur at any time during the development. An assessment should be carried out to help making the monitor policy. Change and supporting are two important factors and need to be assessed. There are many methods for assessing these two factors. For illustration purpose, a simple weighted rating scheme is proposed. Rating scheme is a well-known analysis technique frequently used in the operations management[3,4]. In summary, software development is a dynamic process and change is bound to occur. Thus the change and importance factors need to be quantified before development.

Establishing Management Policy

To establish a set of ground rules from which to monitor, management is encouraged to discuss the project objectives and requirements with all team members at an early stage of software development. In addition to selecting appropriate metrics and standards to measure expenditures, productivity and performance will also be discussed. As a result, all team members have a better understanding of the overall development process, their roles and commitments. In other words, the whole development team will move towards the same goal.

Monitoring Mechanism

A monitoring mechanism need to be established before for collecting all useful information and measuring variances from the work plans. This system should span the software development process from requirements analysis to integration testing. In addition, the basic intention of the system is to support two types of progress monitoring[5]:

- **Milestone Monitoring**
Upon the completion of a particular stage in the development cycle, several milestones had been established to evaluate the progress.
- **Continuous Point Monitoring**
Continuous Point monitoring is similar to milestone monitoring, but can take place at any time during a stage, and utilize whatever data happen to be available. The basic difference between these two mechanisms is that continuous point monitoring involves just one metric overtime and the time of monitoring is fixed. For example, outside supporting level is monitored by the point monitoring at the early of coding stage each week and will forfeit the operation upon the completion of this stage.

Effort managers are free to add and collect more information items in their own performance monitoring system. Obviously, software development is a dynamic process and change is bound to occur. Thus all changes have been closely monitored and justified.

Establishing an Historical Database

It is essential to build a database to ensure continuous improvement on the estimation knowledge. Besides, the database is very useful in making the estimation. It helps the project management in considering the reuse of components of the existing software and determining the potential reusability of the software under design.

OPERATION PHASE

The operation phase includes the activities of monitor, evaluation and determination. The major management tasks in this phase are described below:

Performance Monitoring

Using metrics and standards can keep track of the project progress of all activities in terms of time, size and quality. The performance monitoring mechanism operates continuously in the development process until the products finally delivered to the customer. To get a better monitoring, the authors believe that informal monitoring mechanism also play an important role and the following points are worth to be highlighted.

- Measure the performance of a group rather than an individual.
- Formal and Informal meetings with different teams are worth holding during development cycle.
- Observe the overall performance of an individual in informal basis
- Pay more attention to the schedule of absence from work.
- Pay more attention to the unpaid overtime of an individual.

Based on written report, formal/informal meeting, and observation, the monitoring system can collect different kind of information relating to the actual progress of the development. In summary, the monitoring system assists the management to spot symptoms of problems during the development process.

Evaluating Collected Information and the Monitoring Metrics

The major idea of this stage is to assist the project managers to identify what problem really happened in the project based on all collected information only. In other words, all collected project progress data from the

performance monitoring system must be evaluated and represented in a standard format. Thus it can be interpreted consistently by project management.

Before evaluation, the project managers have to review the evaluation metrics in order to make sure all measurement tools are good to identify or state the problem correctly. In other words, all collected information will be useful and truly reflect the real problem existing in the development process. If the metrics can not truly reflect the real problem, the project managers need to adjust them and re-evaluate all information again.

Evaluating project progress can be regarded as a validation process to ensure that project progress information reflects the actual project situation. Based on all collected information, the expected progress in deliverables against actual progress will be examined and provides the "rough" boundary of the problem. For example, schedule evaluation aims at comparing the actual expenditure of time with the planned schedule. After evaluation, project managers understand whether the project will be finished with an unacceptable schedule. To standardize the evaluation process, some internal standards or international performance standards can be referred.

It should be reminded that some symptoms may be common to different problems or one symptom may be caused by several problems. For example, the symptom of overtime working may be caused by poor communication, poor job allocation, insufficient support, etc. The most important issue is that the project managers need to know problem has happened so that remedies can be devised.

In summary, during the software development cycle, the metrics and standards should be reviewed from time to time. Once some of them are ineffective or inappropriate, they should be removed from the effort management system and may be replaced as necessary.

Determination

Obviously, decisions are essentially responses to problems occurring in software development. Once the evaluation metrics are found appropriate and the problem has been identified, the project managers must determine the boundaries of an acceptable solution and identify all possible courses of action. Considering all updated project information from the monitoring system and other relating factors, project manager needs to determine whether the project has deviated from the plan.

To achieve the goal, a two-step mechanism is adopted. The first step is to re-arrange the project plan based on the updated project information. The project schedule and size plan re-construction will eliminate the self-complementary effect. For example, the early completion of one task and the overrunning of another task may complement each other and make it difficult for the project managers to be aware of the problem.

Having a new project schedule, the second step is to consider all factors relating to the event in order to predict the future result. For example, lower productivity level may be accepted in the early stages of the software development cycle because team members have to spend time to understand the complex application. The project managers need to consider all possible factors before making decision. Thus in each decision making, the project managers explicitly or implicitly make assumptions about the future. If the event will cause the project to deviate from the target time or budget, some corrective actions must be taken.

In general, Acquisition Phase would form part of the activity during feasibility/requirement stage of a "waterfall" software development life cycle model, while operation phase would certainly be operated continuously in the development process until the products finally delivered to the customer. The estimate at acquisition phase may be a "rough" estimate based on limited information only and will be further refined during the development cycle if necessary.

A CASE STUDY

An empirical case study shows how the SFM was applied in real life software development. The project was requested by a large jockey club and developed by an international software consulting organization. This project was designed to replace the existing telephone betting system on a proprietary hardware

platform. The client requested to develop this system on a Personal Digital Assistant (PDA) computer and several advance wireless transmission technologies would be included. The application is a new system and consists of two parts, namely Off-line system and On-line system.

Off-line system

The off-line system allows the user to store details of various bets, the user then must connect the customer terminal (PDA device) to the central Telebet system (betting via tele-communication process) to transmit the details of the betting to the Telebet facility. This process should go through a built-in wireless modem from any suitable telephone hookup worldwide.

On-line system

The objective of this system is to monitor and control the betting process. For example, the system allows users to query the state of the Telebet account balance or to display the results and dividends from completed race. Besides, this function also allows the user to withdraw funds from the Telebet account and transfer the funds directly to the user's bank account. Due to various major elements of uncertainty in this project, a decision was reached to apply SMF. The company derived a SMF and it was expected to provide the management with more information and a reliable control for each phase of software development.

Finally, the project completed and delivered on schedule. The software project is eventually successful. It is not because there are no problems during the development cycle but because the problems are found and overcome before causing serious deviation. One of the major tasks for SMF is the prediction of development efforts and monitors the progress of the development. Once a deviation has been found, project manager need to determine the boundary of the problem. This arrangement ensures that the software project can be developed according to plan and to be ready exactly on time.

CONCLUSION

Experience has shown that most project managers seldom enjoy the luxury of excess resources and spare time in doing the effort estimates. Moreover, the resource constraints and different kinds of uncertainty cause a lot of difficulties in project effort estimating and planning. The success of a software project relies very much on a good management and control system, which allows the development to satisfy the project objectives. However, the development of software project involves many factors, which are hard to manage in this ever changing world. Some examples of these factors are human relationships, technological changes and frequency of requests for changes by users.

As mentioned earlier, software project development has a dynamic nature. Therefore, changes are bound to occur. This paper presents a monitoring framework to the project management in order to tackle the problems from changing software development. Generally, the SMF is developed to let project managers monitor the development process and make the project development running on schedule. In order to use this framework successfully in practice, the following are important requirements:

- close working with team members
- understanding the needs of the users
- commitment of management before starting
- close monitoring of the development progress
- using different points of view to do evaluation at different development stages
- evaluation will be redone while change appears

The authors believe that the SMF and above requirements are beneficial in the monitoring and controlling the development effort of software projects in general. The use of the framework will likely result in improved software development for those project managers who choose to apply it.

ACKNOWLEDGEMENT

The authors would like to thank Dr. Chuk Yau and Mr. T.P. Cheng for their encouragement and support during the course of this research.

REFERENCES

1. Phan, D., Vogel, D. and Nunamaker, J., "The Search for Perfect Project Management", Computerworld, September 1988
2. Heemstra, F J, "Software Cost Estimation", Information and Software Technology, Vol. 34 No 10 October 1992
3. Goodman, P., "Practical Implementation of Software Metrics", McGraw-Hill, 1993
4. Ashley, N., "Measurement as a Powerful Software Management Tool", McGraw-Hill, 1995
5. Barbara, A. K. and Walker, J. G., "A Quantitative Approach to Monitoring Software Development", Software Engineering Journal, January 1989
6. Heemstra, Kusters and Genuchten, "Selections of Software Cost Estimation Models", Report TUE/BDK University of Technology Eindhoven 1989
7. "Betting Rules", The Royal Hong Kong Jockey Club, September 1990

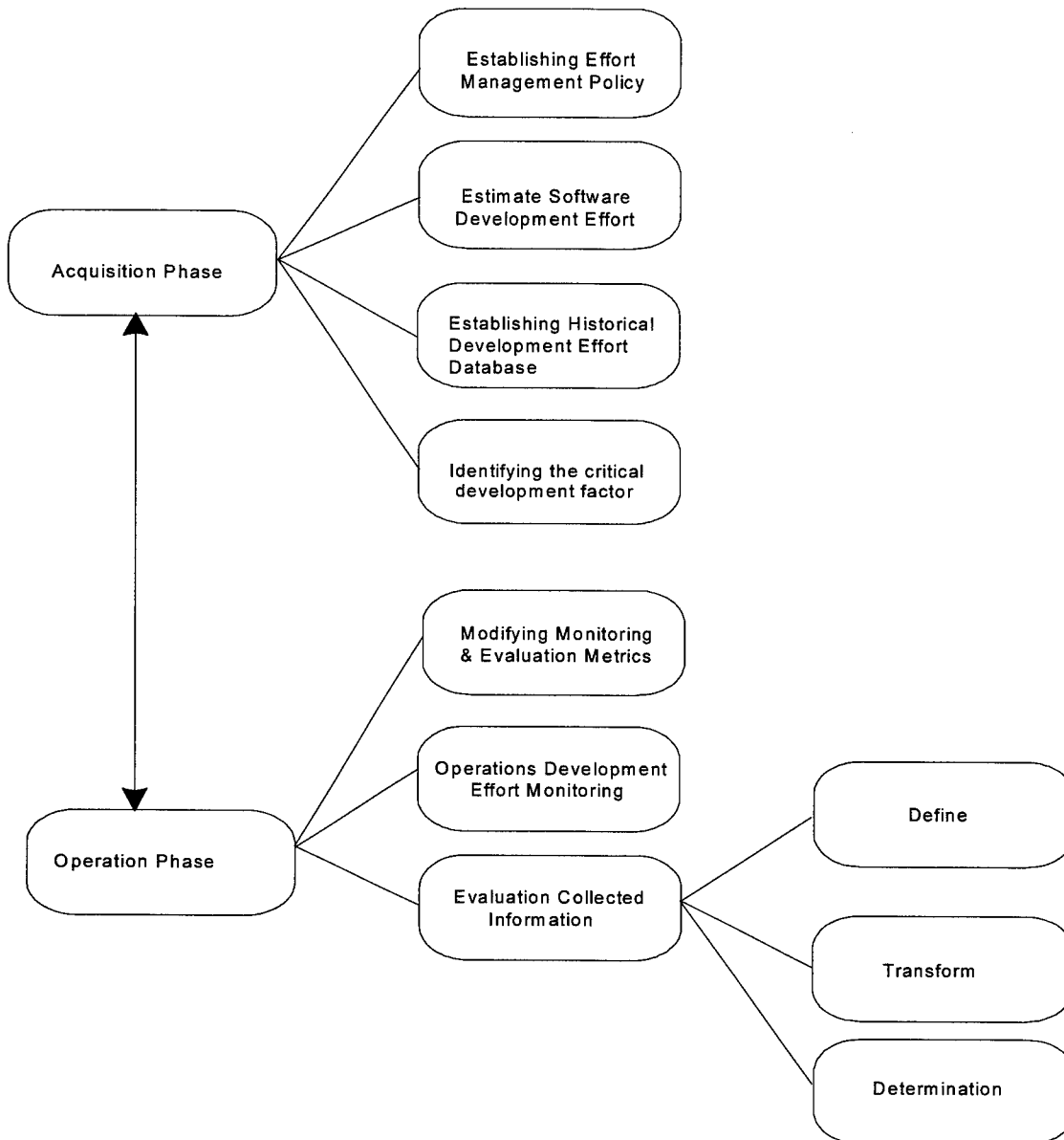


Fig. 1. Software Monitoring Framework.

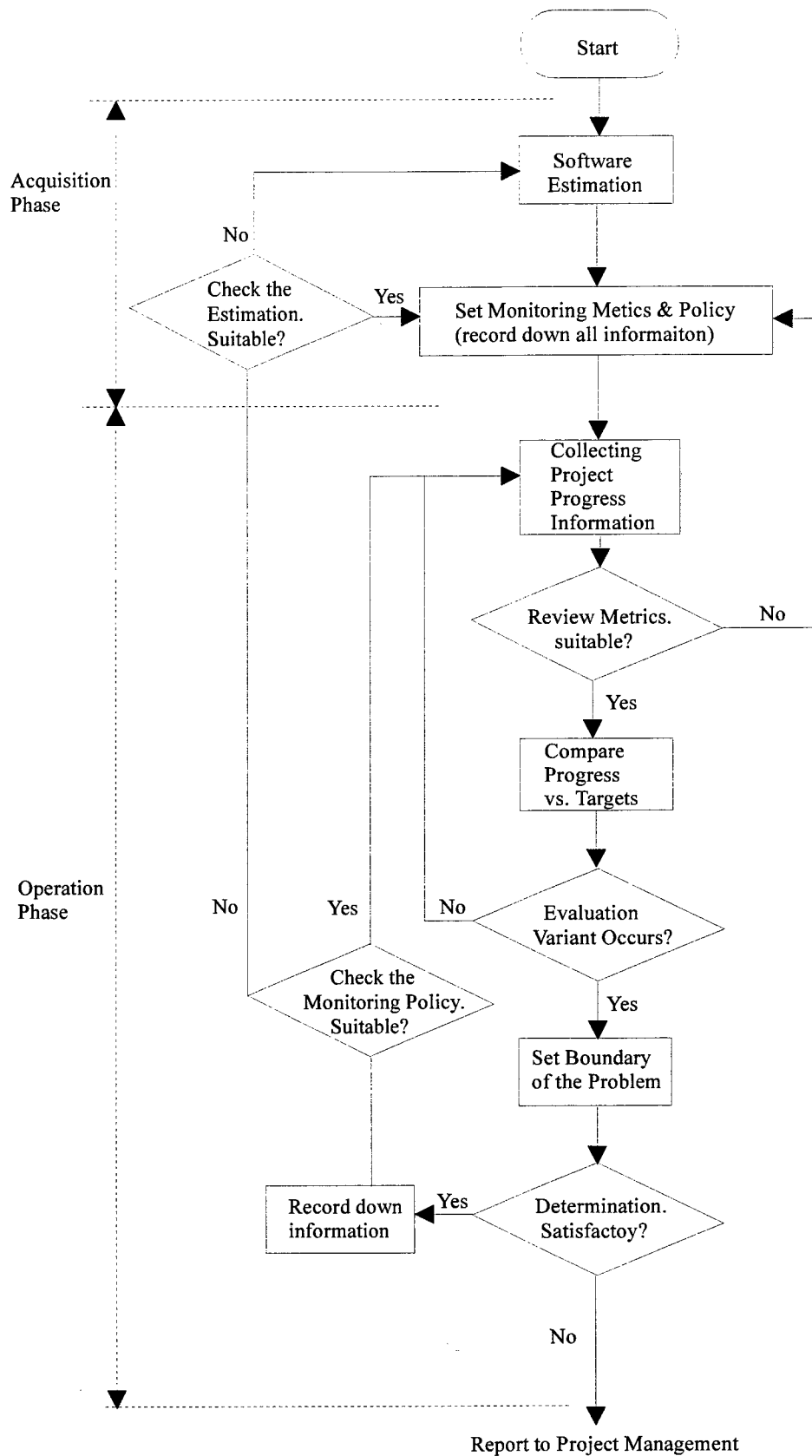


Fig. 2. The Monitoring Work Flow.

Redefining the Web: Toward the Creation of Large-Scale Distributed Applications

Guy M. Nicoletti

Engineering Department, University of Pittsburgh at Greensburg
1150 Mt. Pleasant Road, Greensburg, Pennsylvania 15601

ABSTRACT

The objective of this paper is to illustrate how object-oriented programming is used to generate Distributed Object Programming (DOP): a novel framework for developing multi-user distributed applications that exploit WWW (*World Wide Web*) infrastructure. The paper is organized as follows. Part I briefly introduces notions of HTML, CGI, CORBA, and Java. Part II describes the development process of a multi-user distributed construct in terms of CORBA and WWW-based Java applets. Part III presents an overview of an operational structure of the system developed in Part II. A brief analysis of the (DOP), related issues, conclusions and recommendations are found in Part IV.

INTRODUCTION

In general, implementations of multi-user distributed applications tend to be built using the WWW based HTML (*Hypertext Markup Language*) [1], and the CGI (*Common Gateway Interface*) [2]. HTML is the formatting language that indicates to the browser how to display the contents of WWW documents. CGI, on the other hand, is a mechanism by which WWW servers execute those special-purpose CGI programs that handle specific requests from a WWW browser. In the CGI system, a user interface is established using the HTML FORM tag. The FORM has two main functions: (1) it makes the browser display a form into which the user enters data and (2) specifies the CGI program that will implement the form's operation. Fig. 1 illustrates these functions.

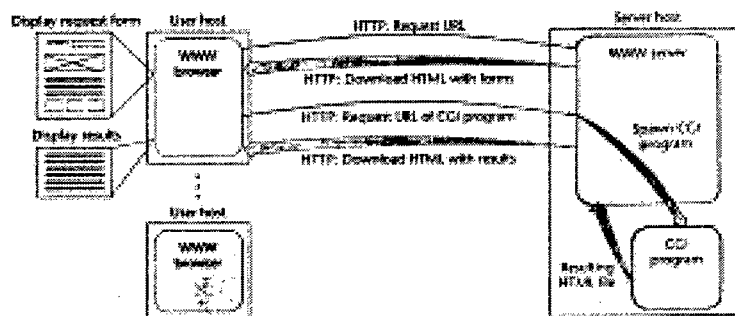


Fig. 1. The CGI System.

With HTML forms and CGI it is possible to (1) add simple user interactivity to WWW sites and (2) gain direct access to a distributed application without the need of a client software on the user's host computer. Thus, essentially *the user's interaction with a distributed application is based on the submission of a form for each request, and each request is served by a separate CGI program expanded as needed by the WWW server*. On the other hand, the (1) limited layout control of HTML forms and (2) lack of continuous execution state on either the server or client side, makes the CGI approach inadequate when integrated applications such as automated process control and distributed simulation are considered. Such applications cannot be effectively served only by a set of disjoint request forms.

The CORBA System

The Common Object Request Broker (CORBA) [3], specifies a system which provides interoperability between objects in a heterogeneous, distributed environment also transparent to the user (programmer). Its

design is based on OMG Object Model. OMG defines common object semantics for specifying the externally visible characteristics of objects in a standard and implementation-independent way. CORBA operates as follows: *clients* request services from *objects* (servers) through a well-defined interface. This interface is specified in OMG IDL (Interface Definition Language). A client accesses an object by issuing a *request* to the object. The request is an event, and it carries information including an operation, the *object reference* of the service provider, and actual parameters (if any). The object reference is an object name that defines an object reliability.

Basic Mechanism of issuing a request.

Fig. 2 shows the main components of the ORB architecture and their interconnections. The central component of CORBA is the *Object Request Broker* (ORB). It encompasses all of the communication infrastructure necessary to (1) identify and locate objects, (2) handle connection management, and (3) deliver data.

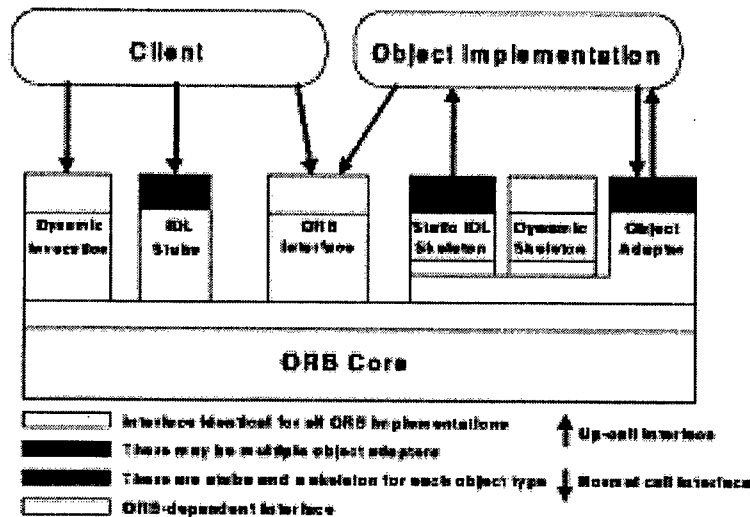


Fig. 2. The ORB Architecture.

The basic functionality provided by the ORB consists of passing client requests to the object implementations on which they are invoked. To make a request, a client can communicate with the ORB Core through the IDL *stub* or through the Dynamic Invocation Interface (DII). The stub represents the mapping between the language of implementation of the client (C, C++, Java, and others) and the ORB core, provided that the implementation of the ORB supports this mapping. The ORB core then transfers the request to the object implementation which receives the request as an up-call through either an IDL skeleton, or a dynamic skeleton.

Architectural Components

A standard component of the CORBA architecture is the Implementation Repository, a database which facilitates access of the OA to information about an object's location and operating environment.

Communication

Communication between the object implementation and the ORB core is effected by the *Object Adapter* (OA). It handles services such as (1) generation and interpretation of object frames, (2) method invocation, (3) security interactions, (4) object and implementation activation and deactivation, (5) mapping references corresponding to object implementations and registration of implementations.

OMG

OMG specifies four policies in which the OA may handle object implementation activation: *Shared Server Policy*, in which multiple objects may be implemented in the same program; *Unshared Server Policy*; *Server-per-Method Policy*, in which a server is started each time a request is received; and *Persistent Server Policy*.

Interfaces

Interfaces to the objects can be specified in two ways: (1) in OMG IDL, (2) as an addition to the *Interface Repository*, another component of the architecture. The Interface Repository also contains information about types of parameters, certain debugging information, etc. Interface to objects of unknown definition at compile time is established by a client's request via the *Dynamic Invocation Interface*.

Server

A server side analogue to DII is the Dynamic Skeleton Interface (DSI). DSI is: a way to deliver requests from the ORB to an implemented object that does not have compile-time knowledge of the object to be implemented, an interactive software development facilitator based on interpreters and debuggers; and a provider of inter-ORB interoperability.

Interoperability

There are many different ORB products currently available. Furthermore, there are distributed and/or client/server systems which are not CORBA-compliant. In order to provide needed interoperability, OMG has formulated the ORB interoperability architecture [4].

The JAVA System

The Java system is a new programming structure which has gained world-wide rapid acceptance due primarily to its *object orientation*; *platform independence*; and *networking capabilities*. The pillars of Java structure are the innovative programming language; and *Java Development Kit (JDK)*. These two components make programming convenient, efficient, and effective. The JDK includes tools for building GUIs, networking, and implementing *applets*, *small Java programs that add dynamic content to Web pages*. Object-Oriented Programming (OOP) makes software easier to build, maintain, modify, and reuse. With OOP one programs by building *software objects*. Because of *platform independence*, Java programs are compiled into *Java bytecode* that runs on any computer with a *Java interpreter* or *Java Virtual Machine (JVM)* providing a run-time environment for executing Java programs. The bytecode is architecture and operating system independent. Two main features: Java's bytecode and the HTML APPLET tag are especially attractive and useful for distributed applications. By extending the WWW browser to incorporate a Java runtime, then the APPLET tag provides the necessary information for a WWW browser to find, download to the user's host computer, and execute a Java applet. Thus, the *WWW infrastructure-Java applet* embedding simplifies deploying of user client components of distributed applications [5].

THE WWW/JAVA/CORBA SYSTEM

As implied earlier, multi-user distributed applications are driven by their interaction with users, who invoke operations implementable by server software residing on remote host computers. Obviously then WWW, the Java language environment, and CORBA architecture are indeed complementary software technologies, which, when combined provide a new and powerful construct for developing and deploying such multi-user distributed systems. Specifically, one can conceive of building a "user-friendly" client software as WWW-downloadable Java applets, which employ CORBA's ORB and IDL to interact with remote server software. This section describes the structure of such software system. This is illustrated in Fig. 3.

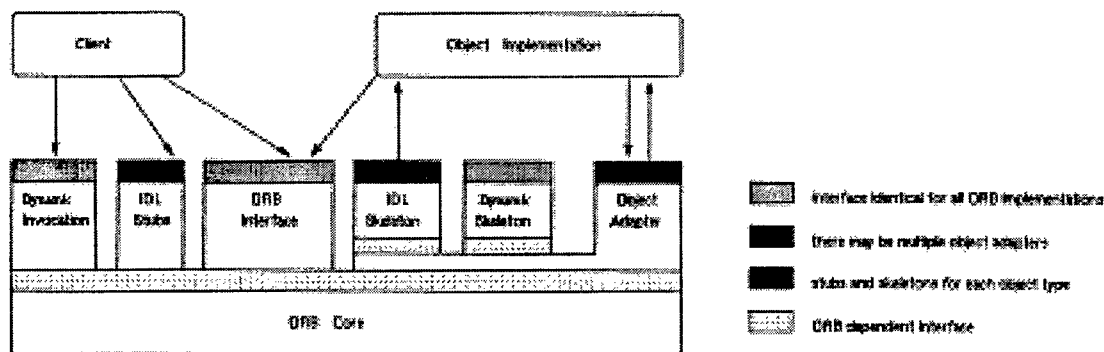


Fig. 3. CORBA'S ORB-IDL Interaction.

This structure may be modularized as follows:

Modulo 1: Allocation of tasks

Definition and allocation of tasks between user's host and server site must be established. Typically, the user's host does such tasks as *user interfacing process, user input error checking and multithreading for user-GUI interaction*; while *shared resources, related software to control them, and remote objects* are allocated to the server.

Modulo 2: Remote Object Interface

Here, interfaces of the remote objects are defined in terms of the operations the client software can invoke on the remote objects. These definitions are organized in an IDL file. Major server-side as CORBA interface may be so modeled.

Modulo 3: Implementation of Remote Objects

In order to automatically generate software for the remote object's server program, remote objects are implemented on an IDL compiler. This is a skeleton code that includes: type definitions in the IDL; code for intercommunication between incoming requests and remote objects; and the empty bodies of the methods. Next is the implementation of the source code of these methods to provide the remote object's capabilities. Finally, use the remote object's implementation class to build a server program. Such a server program may be built in C++ CORBA Basic Object Adapter (BOA) in shared server mode. OrbixWeb may be used to accomplish this task [6].

Generate Remote Object Proxies

Implement the Java applet that invokes method on the remote object. This is accomplished with an IDL-to-Java compiler that generates the client-side proxy class that represents the remote object class. A proxy class for a remote object contains the information of the address of the remote object. Thus, when the client invokes a method on a proxy object, the proxy relays that invocation to the real remote object via the ORB by OrbixWeb products.

Obtain Remote Object References

Here, integration of method invocations on the remote object into the Java client applet is established. The client obtains a reference to an instance of the remote object's proxy class. Either OrbixWeb or CORBA accomplishes this task. With CORBA, however, one can also define a more general-purpose mechanism in the standard Naming Service Interface (NSI) thus allowing remote object servers to register their remote objects by name [7]. In addition, applets could use this mechanism to look up any number of other remote object references by name.

Invoking Remote Methods

The next step is for the client to invoke methods on the object. This object is regarded as the real version control object. All the transport-level messages and data handling needed to convey the invocations of the remote method between client and server are implemented with the ORB library software, and the software generated by the IDL compilers.

System Configuration

To configure the system one must: install the remote server program on the server host; make its remote objects' references available through one of the mechanisms adopted in modulo F above; install the compiled Java class files that constitute each client applet on the WWW server host so they can be downloaded by a browser; and for each applet, install an HTML file on the WWW server. Each such file must include the HTML APPLET tag that specifies the applet's location on the WWW server. The system is then configured to operate as follows.

1. The client software is deployed automatically (since it is implemented as applets).
2. As a user opens a URL for an HTML file containing an APPLET tag, the browser downloads the applet's executable JAVA class files, then runs the applet under control of its own Java runtime.
3. To request an operation on a remote object, the applet invokes the operation on the object's proxy; the proxy relays the request to the real object using CORBA protocols.

Essentially, the described system, intends to address the difficulties associated with network programming of Distributed Applications (DAs) which, often consist of several network communicating programs written in different programming languages and running on different operating systems. Thus, this system must provide a software environment that allows the developer to build DAs that interact as though they were implemented in a single programming language on one computer. Such software environment is to be provided via an implementation of CORBA from the Object Management Group (OMG) that maps CORBA functionality to the flexible Java programming language. This is a huge and complex undertaking. One such system, however, is to be found in the proprietary IONA's OrbixWeb system in which applications are fully Internet-enabled [7]. Within OrbixWeb, CORBA defines a framework for developing interactive object-oriented, distributed applications; facilitates design of a distributed application as a set of cooperating objects; allows reuse of existing objects in new applications; and maps applications of its distributed programming into the Java environment which is described next.

THE WWW – JAVA -CORBA ENVIRONMENT

Essentially, the implementation of an online Java-CORBA Distributed Application in a DOP environment consists of the steps described below [8].

- Step 1. *Develop CORBA programs in Java using WWW site environment.* Programming employs the CORBA Interface Definition Language (IDL) and standard OMG mapping from IDL to Java using proprietary ORB. The IDL to Java mapping comprises mapping for (1) basic data types, (2) modules, (3) interfaces, (4) constructed types, (5) strings, (6) sequences, and (7) arrays.
- Step 2. *Create and manipulate Java Applets.* define the web *Implementation Repository*, provide the server, write client Applets, add Applets to HTML files, and add Web Client Functionality.
- Step 3. *Implement IDL Interfaces.* The programming steps conducive to IDL interfaces comprise (1) Interfacing application objects, (2) Compiling the IDL interfaces, and (3) Implementing the interfaces. This last step leads to the development of the server application and of the client application, followed by Registration-Activation, and Execution Trace.
- Step 4. *Establish ORB Interoperability.* ORB interoperability enables the client of one ORB to invoke operations on an object in a different ORB via an agreed protocol. This is accomplished via the two OMG specified standard protocols, namely the General Inter-ORB Protocol (GIOP) - which defines the on-the-wire data representation and message format; and the Internet Inter-ORB Protocol (IIOP). The IIOP is an OMG-defined specialization of GIOP that uses TCP/IP as the transporter layer.
- Step 5. *Run ORBWeb Clients.* The following requirements must be fulfilled: (1) obtain access to the Java code for an application, (2) make the code available to the Java bytecode interpreter, and (3) run the interpreter on the class containing the *main()* method for the application.
- Step 6. *Registration and Activation of Servers.* This is done via the Implementation Repository (a database of server information). The Implementation Repository is implemented in the OrbixWeb *daemon* -- a Java implementation of the program interface. The Implementation Repository maintains registration information about servers and controls their activation.
- Step 7. *Establish Dynamic Skeleton Interface.* The Dynamic Skeleton Interface (DSI) is the server-side equivalent of the Dynamic Invocation Interface (DII). The DII allows a client to call operations on IDL interfaces that were unknown at the time of the client compilation. Thus, the DSI allows a server to receive an operation or attribute invocation on any object, even one with an IDL interface unknown at compile time.
- Step 8. *Establish Interface Repository.* The Interface Repository (IFR) is the OrbixWeb component that provides persistent storage of IDL interfaces, modules and other IDL types. A program can browse through or list the contents of the Interface Repository. A client can also add and remove definitions from the Interface Repository.
- Step 9. *Generate Smart Proxies.* Smart proxies allow optimization of client interaction with remote services. The IDL compiler automatically generates proxy classes for IDL interfaces. When a proxy

receives an invocation, it packages the invocation for transmission to the target object in another address space on the same host, or a different host.

Step 10. *Provide for Loaders.* *Loaders* are designed to support persistent objects – long-lived objects stored on disk in the file system or in a database. When an invocation arrives at a server process, OrbixWeb searches for the target object in the internal object table for the process. By default, if the object is not found, OrbixWeb returns an exception to the caller. If one or more *loader* objects are installed in the process, OrbixWeb informs the loader about the object fault and allows it to load the target object and resume invocation transparent to the caller.

CONCLUSION

Based on the author's experience with the OrbixWeb Demo Manager, the Applet authentication remains a partial problem due to browser-imposed restrictions on file-based URLs. These restrictions depend on the type of browser being adopted. There are other issues; however, these are of minor consequence to the operation of the environment. Such issues are being systematically eliminated with sequential software upgrading. A major advantage of this WWW-Java-CORBA system is that it matches the Java Remote Method Invocation (RMI). With OrbixWeb, the IDL module, the Remote Object Server, and the Class Software Objects are arranged in an integrated system. Integration is based on the structure of the OrbixWeb "Daemon". IDL interface to the OrbixWeb is such that (1) the daemon is responsible for activating servers (if an appropriate server is not already running) and dispatching operation requests. (2) The daemon is involved, if at all, only with the first operation request from a client – it is not involved with subsequent requests. (3) Two OrbixWeb daemon executables are available (the Java Daemon). (4) The Daemon is also responsible to manage the Implementation Repository, to search for an appropriate server via the locator and to manage the following configuration files used by the default locator: the *server location file*, and the *host groups definition file*. In addition to its deployability and platform independency, this WWW-based Java/CORBA approach offers the following advantages over both the CGI and the CORBA approach.

- *The Internet ORB* -- establishes tight integration a proprietary firewall for Internet Protocol communication, ensuring that applications are fully Internet-enabled.
- *The Java Server* -- The system includes a Java version of the server activation component and associated utilities with which it is possible to launch Java CORBA server components as new threads or as multi-thread deployment.
- *IDL/Java Mapping* -- Facilitates full support for the standard OMG IDL to Java interface including the Java ORB portability interfaces.
- *Fragmentation* -- Provides Internet Inter-ORB Interoperability Protocol which allow one to send IIOP messages as fragments. This increase parallelism, improves the overall dispatch speed for very large messages, and fosters lower memory consumption.
- *CORBA Services* -- In addition to the set of standard CORBA interfaces, this module includes such services as Naming Service, Transaction Service, and the Event Service, among others. CORBAService offers a level service standardization not present as yet in GUI Java.

In conclusion, this paper has presented a more flexible and more powerful distributed object technology consisting of interfaced CORBA, Java RMI and Java applets deployed via the standard and widespread WWW. This technology provides, among others, a sophisticated set of server-based capabilities remotely accessible by users.

REFERENCES

1. T. Berners-Lee, D. Connolly, 1995. "Hypertext Markup Language – 2.0", Internet RFC 1866, Nov.
2. <http://www.w3.org/pub/WWW/CGI/Overview.html>.
3. <http://www.infosys.tuwien.ac.at/Research/Corba/OMG/arch2.htm>
4. <http://www.expersoft.com/Resources/DistTech/tutorial.htm>
5. E. Yourdon, 1998. "Java, the Web, and Software Development", *Computer*, 29(8), 25-30.
6. Iona Technologies PLC: <http://www.iona.com>
7. Iona Technologies PLC: <http://www.iona.com>. "OrbixWeb Programmer's Guide", Sept. 1998.
8. Iona Technologies PLC: <http://www.iona.com>. "OrbixWeb Programmer's Reference", Sept. 1998.

How Can We Form/Expand Conceptions in Workers' Minds According to Their Individualities?

Kumiko Ishino

Konan University, The Department of Science,
308, Residence RS, 3369-1, Komanyu-Machi, Utsunomiya-City
320-0065, Tochigi-Pref, Kobe, Hyogo, Japan
Email: kumiko.i.m@ma4.justnet.ne.jp

ABSTRACT

Every person has individuality. We need various kinds of supporting aids designed according to our characteristics. Workers also need supporting aids designed according to their individuality. When generating supporting aids, we should consider the different effects of many types. We should define as attributes of the individuality not only states of logic but should also include comprehension states, thought processes and individual abilities. This paper tries to express the individuality of workers, and defines the formation and expansion of physics concepts. As a result, I will propose methods to expand and form conceptions for the worker according to their individualities by providing a cross-reference support system.

INTRODUCTION

Caring for the individualities of workers requires supporting methods designed for them. We showed that human beings have different types of brains [1]. This means we have individuality, and we need suitable support. This idea is also supported by Denis who showed various thought processes [2]. We can also refer to many types of thinking in history of science [3,4]. Thus, demonstrating the necessity of supporting methods to be designed, requires it to be adaptable to workers' individualities.

Experimental research to define a persons' individuality is often conducted without using support processes. Experimentation accompanies mutual effects. When research is conducted into processes of cognition, it is sometimes passive, as we don't want to affect the process output. However this process involves mutual effects. The substance of support is an example of mutual effect. Therefore, we must conduct research to understand and provide background knowledge to demonstrate the relationship between workers' abilities and backgrounds and their relations with the structure of background aids. Concerning this, Monk spoke of the importance of referring to and using the eyes of a physical scientist to model or interpret data to support learners [5]. As an addition, I propose to use scientific methods, especially discovery processes, to study human science which is extremely complicated. We simplify the problem by referring to discovery processes of science. In physics for example, discovery begins by creating an initial hypothesis and then verifying. I pay particular attention to the importance of hypothesis creation.

In this paper I use as an example, a hypothesis about formation/expansion of a concept and then show how workers' comprehension of statements about universal gravity, affects learning, forming and understanding. I have created methods to support workers as they try to comprehend statements using an intelligent support system called '*Contact*'. Verification of the hypothesis and the effects of the method are future goals.

Regarding this research, I have shown the process of scientific thinking and learners' abilities and proposed educational methods about classical electromagnetism [6]. In 1996, Ishino, Sugai and Mizoguchi proposed an intelligent educational system called '*Galileo*' which supports scientific thinking about movement of an object [7,8]. This paper considers comprehension of a concept in detail, especially about a state '*Galileo*' recognized as "unrecognized". I propose methods to support *Galileo*. In 1997, I proposed aids to help workers' thinking, designed according to their comprehension states, abilities, intentions, working environment and states of their jobs [9]. This paper proposes support methods to form/expand conceptions in workers' minds in detail, according to their intentions and recognition states about concepts and abilities.

There is much research conducted to define development of human recognition states. However, since, humans often respond out of curiosity, intentions of researchers don't often agree with those of their human subjects, so useful information is not obtained [10]. I propose support methods designed to recognize a worker's state by using simple dialogues. We can control the scene required using visual aids.

DEFINITIONS OF FORMATION AND EXPANSION OF CONCEPTIONS AND CONSIDERATION OF COMPREHENSION STATES AND ABILITIES

First, I create a definition about formation and expansion of a conception in workers' minds as follows:

"Physics concepts are formed by understanding a condition in which the value of an unformed concept changes, which is accompanied by a mental image. The concept is defined by providing a name. The concept is expanded by discovering differences with other concepts, or by defining causal relationships with other concepts in a particular scene."

This formation and expansion of conceptions require the following:

- Ability to make a problem concrete
- Ability to operate tools
- Ability to abstract a causal relationship
- Ability to mentally visualize a scene
- Ability to mentally verbalize a scene
- Ability to discover a condition value change
- Ability to define a causal relationship
- Ability to conclude information
- Ability to retain information in the mind
- Ability to name a concept

These individual abilities should be recognized and aids must be tailored to fit the workers' needs.

Suppose a worker is in a situation in which universal gravity is effective. At this time, let us present a scene such as (1) to her/him to confirm if s/he recognizes the concept of universal gravity by asking (2)

- Contact: "A piece of freight falls from a shelf now." 1.
Contact: "Will universal gravity affect the piece of freight?" 2.

If the worker answers "No, it won't.", then we realize that s/he doesn't recognize universal gravity for one or more of the following reasons:

- (a) S/he may not have a conception of universal gravity.
- (b) S/he may have a conception of universal gravity but not comprehend the word "universal gravity".
- (c) S/he may not recognize the earth as an object.
- (d) S/he may not recognize the piece of freight as an object.
- (e) S/he may not have a conception of freight.

Concerning aids, we must consider a worker's comprehension state. When creating methods and aids to form and expand concepts, we should recognize the states of a worker from various viewpoints. Concerning this, Domenech, Casasus and Domenech proposed a classification of the most current definition of mass from its physical representation and topic area [12]. They showed the concept as being constructed from different viewpoints. So methods to organize learning aids must also consider different viewpoints.

METHODS TO SUPPORT COMPREHENSION OF A STATEMENT

In this section, I propose methods to help workers comprehend a statement such as (1) above.


Method to allow a worker to recognize the meaning of the name of a concrete object

Let us present a question (3) to the worker to know the reason s/he doesn't recognize universal gravity. Suppose s/he answers question (3) that no, s/he doesn't, s/he may not have the concept, s/he may have the

conception but s/he doesn't know the word "freight", or s/he may know the word but may not be able to visualize it. Furthermore,, her/his ability to retain information may be low. The worker can remember or form a conception by pointing out a relation with a more general concept as in (4) and perhaps by showing an image (5). If we knew the comprehension state of the worker in more detail by using previous communication, we may be able to support her/him by using the appropriate form.

Contact: "Do you understand the meaning of 'freight'?" 3.

Contact: "A freight is a container that gets full." 4.

Contact: "A piece of freight looks like  ." 5.

In the case that an object concept isn't clear in a worker's mind, this method clarifies it.

Oral Method to help a worker recognize an object concept as a concrete object

Confirming whether the worker recognizes the freight and the earth as concrete objects can be established by the following question:

Contact: "Select all concrete objects from the following list: (a. Dream, b. Freight, c. Idea, d. Earth)." 6.

If the worker answers "Both b. and c", then s/he may have recognized them as concrete objects but s/he may not have the conception of universal gravity for this scene. Alternatively, the worker may not have recognized (one of) them as the concrete object(s), and may have simply have selected by guessing. If s/he doesn't select (one of) the concrete objects, s/he may hold the concept of universal gravity but simply doesn't recognize them as concrete object(s). At this point we can support the worker as follows:

Contact: "The freight and the earth are the concrete objects." 7.

or

Contact: "The freight is also a concrete object." 7'.

Oral Method to shape a conception

Next, we consider if the worker has the conception of universal gravity by asking question (8):

Contact: "Will any thing affect the freight?" 8.

If the worker answers question (8), "Yes, something will", s/he may not know the words "universal gravity" but s/he may have the concept and s/he may recognize a causal relation between the freight and universal gravity, or s/he may not recognize the conception of universal gravity but may understand a causal relationship between another force and the freight. Question (9) can help to shape the worker's conception.

Contact: "In what direction will it cause the freight to go?" 9.

If the worker responds that it causes the freight to move directly toward the earth, s/he doesn't know the words "universal gravity" but s/he has the concept and s/he recognizes the causal relation between the freight and universal gravity. Then , by a method which allows the worker to name the concept, the worker can equate the words "universal gravity" with the concept.

Oral Method to expand a conception by relating to a different scene

If the worker responds to question (8), "No, nothing will affect the freight", then it is clear that s/he does not recognize the concept of universal gravity in the scene. If the worker can recognize the concept in another scene, I propose an oral method to expand a conception by relating to another scene. First let us provide another typical scene such as (10), and provide previous oral methods to shape the conception:

Contact: "The moon revolves around the earth." 10.

Contact: "Does anything affect the moon's rotation?" 11.

If the worker answers , "Yes, something does", then s/he may not know the word "universal gravity" but s/he may have the concept and may recognize a causal relation between the moon and universal gravity or s/he may not recognize the conception of universal gravity but may recognize a causal relation between another force and the moon. Question (12) helps to shape the worker's conception:

Contact: "In what direction does it cause the moon to go?" 12.

For this, suppose worker A responds that it causes the moon to move in a circle around the earth, and worker B responds that it causes the moon to move directly toward the earth. As the moon is a type of object, the response from worker A means that universal gravity affects an object from the opposite direction in which the object moves. This is incorrect, therefore we should support her/him in the future.

But worker B doesn't know the word "universal gravity" but s/he has the conception and s/he recognizes the causal relation between the moon and universal gravity. Then, by a method to allow a worker to name a conception, the worker can connect the word "universal gravity" with the conception.

Now worker B has the conception of universal gravity in this scene. Therefore we can expand this conception to a previous scene. At this time, I define that low ability as the inability of the worker to abstract a causal relationship from a specific scene to a causal relationship in another scene. If the ability of worker B to abstract a causal relationship in one scene is high, I propose allowing her/him to recognize the concept in the previous scene by applying this causal relation to the previous scene in the form of question (13):

Contact: "Does universal gravity affect the freight's descent toward the earth?" 13.

Contact: "Since the moon is a type of object, this means universal gravity affects an object from the direction in which it is moving."

If the worker's ability to mentally visualize a scene is low, we can provide visual aids to confirm this fact visually. In this case, we should support her/him according to the ability to make a problem concrete and to discover that a condition value changes. By these supports we can expand a concept in one scene to others.

Organization of computer visual aids to allow workers' to form physical conceptions

In a third situation, worker C responds to question (11), "No, nothing does". This worker doesn't recognize the conception of universal gravity. I now discuss the use of aids to form a concept in this type of worker's mind. First, I propose using visual aids to form physical concepts by measuring the concept or by connecting related concepts visually. These are organized according to a worker's individuality using these rules:

- <1> If values of the conception that should be formed in the worker's mind are not 0 on the earth, the system provides a scene on the earth.
- <2> If the system doesn't provide a scene on the earth by reason that the value of the concept is 0 on the earth, the system selects and provides a typical essential scene for the concept.
- <3> If values of the concept are accompanied by objects, the system provides objects by a button labeled "select an object".
- <4> If values of the conception can be measured by a tool(s), the system selects and provides the tool(s) according to the scene in question.
- <5> If the worker recognizes a conception(s) which relate(s) to that which should be formed in her/his mind and which can be measured, the system selects and provides a tool(s) by which the worker can measure the values of the concept according to the scene.
- <6> If the worker has high ability to define a causal relation, the system selects and provides tools by which the worker can measure values of conceptions which relate the conception s/he should form in her/his mind but in which s/he is unable to define the causality.
- <7> If the values of different tools are similar, the system selects a scene(s) where values are different.
- <8> If the worker forms the targeted conception, the system selects and provides a scene(s) where s/he recognize the conception using variables ranging from zero to infinity.
- <9> If the system provides a new scene, the previously selected tools are provided as well.
- <10> If the provided tool(s) don't work in the new scene, the system selects tool(s) that can measure the values of the conception using previously targeted tool(s) and provides it/them in every new scene.

Now let us consider that worker C has the conception "mass of objects" which is measured by weight and can be demonstrated using a balance beam but her/his ability to define a causal relation in a scene is low. Fig.1.(a) is provided as visual support for this type of difficulty to allow recognition of a causal relation

between universal gravity and mass. So, if s/he pushes the button "to the moon", the scene on the moon as in Fig. 1(b) is provided. S/he can also compare the universal gravity which affects objects on the moon with the universal gravity that affects objects on earth which are demonstrated using the spring scale aid. Furthermore, s/he may also weight the mass of objects using the simulated balance beam, wherein the causal relationship between universal gravity and mass can be recognized.

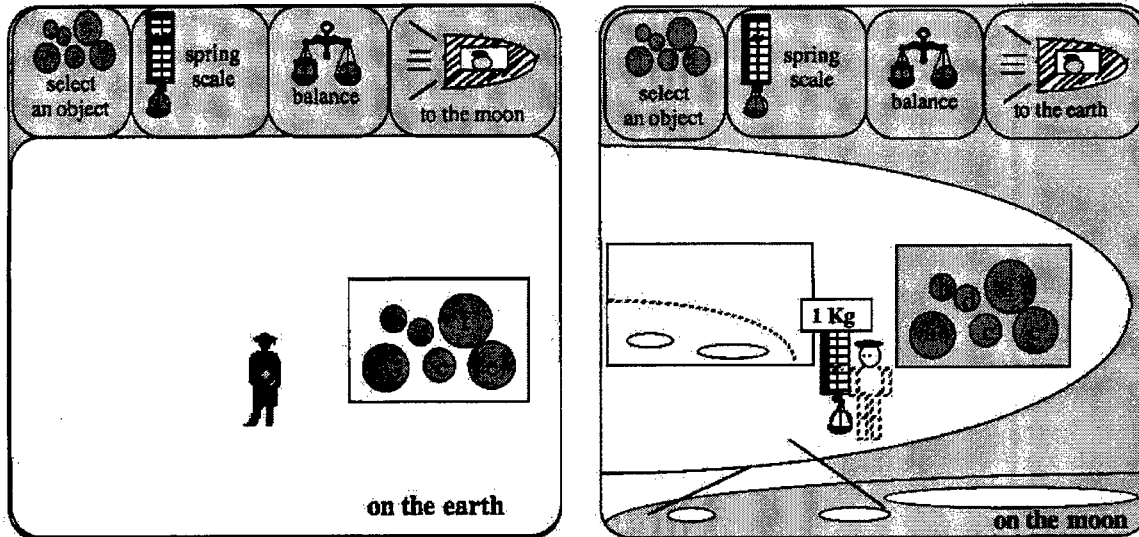


Fig. 1. Computer Visual Aids Contained in 'Contact'.

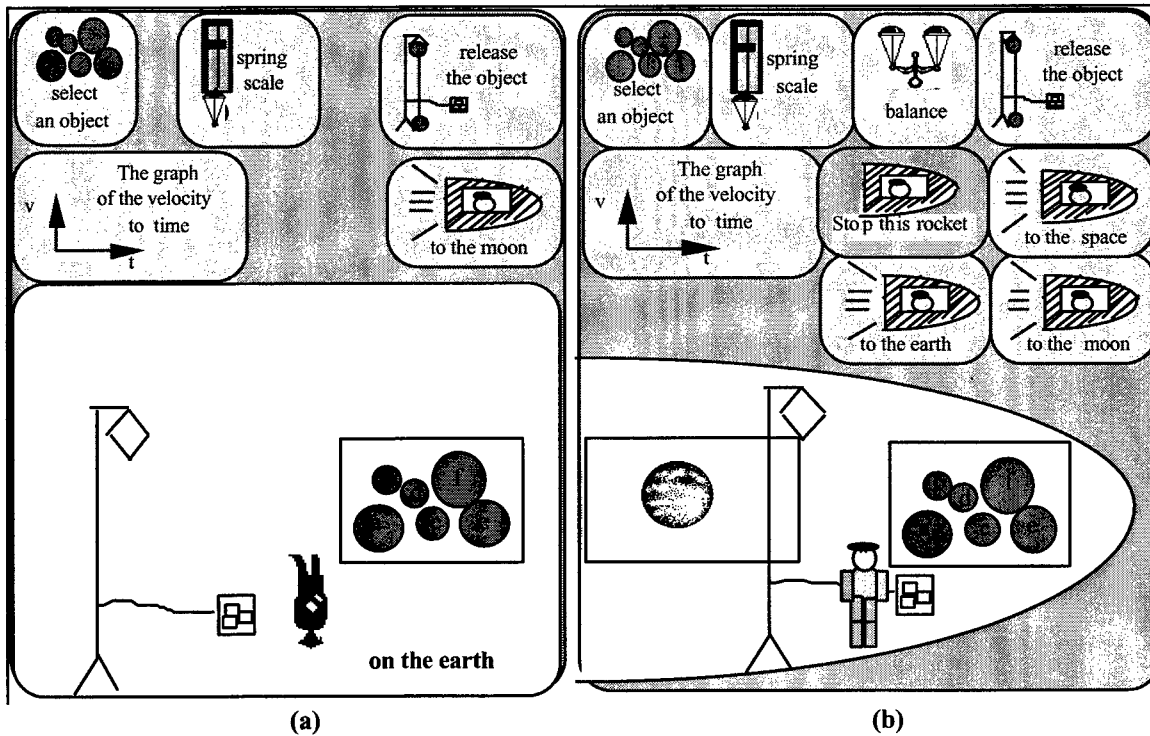


Fig.2.Computer Visual Aids contained in 'Contact'.

However, if the worker has the concept of acceleration at this time, an environment is provided to allow the recognition of the concept of gravitational acceleration and to observe the relation between universal gravity and acceleration (see Fig. 2(a)). Additionally, if the worker has a high ability to define a causal relation, a scene which provides more tools to measure values of conception(s) in which the value of universal gravity changes is provided to allow the worker to discover causal relationships as desired (see

Fig. 2(b)). Since all workers are not able to use the tools and recognize causal relations, we must provide support according to their abilities and/or intentions.

Aids to support operation

If a worker's ability to operate tools is low, s/he may not take action and utilize various causal relation scenes. In such cases, supporting operation should be to guide the worker through the supporting aid system. First if a learner doesn't take action, we can provide instruction (15).

Contact: "You can operate the tools by pushing a button." 15.

If a worker pushes the button "spring scale" but can't measure values of objects, we can suggest:

Contact: "Please select an object and weight it using the spring scale." 16.

Method to create strong motivation

If a worker doesn't take action even after supporting operations have been provided, s/he may not have strong motivation. I propose the following method to create motivation. Suppose we try to allow a worker to form a conception and we want to allow her/him to create motivation. I propose providing an applicable situation about a tool that can measure the concept.

The information would be provided incorrectly using a character in the situation. The character is confusing the conception the worker is trying to form with other conceptions that the worker previously comprehended. Therefore it requires the worker to utilize previously comprehended conceptions for understanding the conception presently being considered. If we provide a correct statement, the worker can only confirm it. But if we provide an incorrect statement, s/he must use her/his knowledge and understanding of conceptions to discover the truth and confirm it by herself/himself. If s/he thinks the statement is correct after confirming it, the program informs her/him that s/he is incorrect. It creates questions in the worker's mind. The following is an example question: "What concept can we measure by the tool?" In this way we provide motivation by appealing to people's natural drive to discover the truth.

Now suppose we try to allow workers to form the concept "universal gravity". We provide false statements about the spring scale which can measure "universal gravity".

Contact: "Doctor A invented a tool 'spring scale". He says the tool can measure mass. Is this true?" 17.

This question applies the conception to a real-life situation. If a worker responds, "Yes, it is", this worker is already confusing the concept of "universal gravity" with the concept of "mass". S/he may not understand the definition of mass, or s/he may not recognize facts. Next we should confirm if the worker understands the definition of mass by question (18):

Contact: "Is the mass of an object changed by moving in space?" 8.

If s/he answers question (18), "Yes, it is", we also must allow her/him to comprehend the concept of mass. If the worker answers question (18), "No, it isn't", s/he understands "mass", but cannot recognize that values of spring scales for an object are changed by moving in space. Therefore I propose allowing the worker to perceive this as follows:

Contact: "Confirm it." 19.

If the worker can't take action, her/his ability to make a problem concrete is low. I propose aid to make the problem concrete as follows:

Contact: "Compare weight values on the moon with weight values on earth by using the spring scale." 20.

If a worker can measure values of objects using a spring scale but s/he doesn't compare values of universal gravity on the earth with those on the moon as a result of not pushing the button "to the moon", we can provide guidance as follows:

Contact: "Let's go to the moon." 21.

After the worker conducts a series of experiments, we can confirm if s/he recognized that the values of the spring scale on the moon are different from those on earth by asking question (22):

Contact: "Are the values equal on both the moon and on the earth?" 22.

If the worker responds, "No, they are not", s/he recognizes these facts and s/he may perceive that the "spring scale" does not measure mass or s/he may have a question about mass changing by moving in space. We can ask her/him about it by using question (23):

Contact: "So, is a 'spring scale' the tool to use to measure mass?" 23.

If the worker responds, "No, it is not", s/he understands "spring scale" can't measure mass. Her/his mind has changed and s/he may have strong motivation to answer unanswered questions such as:

- "What is the value of the spring scale?"
- "What tool can measure the value of mass?"

If a worker responds to question (23), " Yes, it is", her/his ability to conclude or ability to retain information is low or her/his belief in the concept of mass is wavering. We should provide an aid to assist the worker to come to the conclusion in statement (24)

Contact: "Mass of an object doesn't change due to moving in space, but the values of a spring scale do. Therefore a spring scale is not the tool to use to measure mass." 24.

Method to allow worker to discover conditions in which values of a new concept change by actions

After a worker measures the values of the new concept visually, we can confirm if s/he has formed the concept by using question (25):

Contact: "What concept can be measured with the spring scale?" 25.

If the worker responds, "It is a type of force", s/he has recognized the concept. But if s/he can't respond, her/his ability to discover that a condition value changes may be low. For a worker like this, I propose to allow the worker to consider the condition that the value of a new concept changes by actions as follows:

Contact: "How can you get a large value from the spring scale?" 26.

Method to allow a worker to consider meaning of an action

If a worker responds to question (26), "By pulling the spring with my hand", I propose to allow her/him to consider the meaning of the action. Question (27) allows her/him to think about what is measured by pulling a spring scale with her/his hand. If the worker cannot respond to question (26), her/his ability to mentally visualize/verbalize a scene is low. So we can assist her/him by showing the scene visually/verbally.

Contact: "What do you add to the spring by pulling on it?" 27.

If the worker responds, "I add a force", then s/he has formed a new concept of a type of force.

Method to allow a worker to name a concept

After a worker has formed a new concept I propose a method to name the concept because by naming the new concept, it will become solidified. If there is a public name for the concept, we can show it as follows:

Contact: "We call this force 'universal gravity'." 28.

If there isn't a name, we can allow her/him to name it as follows:

Contact: "What do you want to name it ?" 26.

The worker may name it "starforce", and so the concept is solidified in her/his mind. But if the worker balks at giving a name, we know her/his ability to name a concept is low. We can provide assistance as follows:

Contact: "What do you think about the name 'G'?" 30.

Method to allow a worker to discover the difference between a known concept and a new one

If a worker forms a new conception and s/he doesn't understand the differences between a known conception and the new one, I propose a method to allow her/him to discover the difference. Suppose a worker knows the concept of mass and now s/he forms a conception about universal gravity. Computer visual aids such as those in Fig. 1 and Fig. 2 are available. We can confirm if s/he has discovered the difference between these concepts by using question (31):

Contact: "What is the difference between mass and universal gravity?" 31.

If the worker responds that the mass of an object doesn't change but universal gravity for an object changes according to the scene, then s/he understands the difference. But s/he can't respond to the question (31), s/he hasn't discovered the difference. After that, if the worker compares the values of universal gravity with the values of mass on the earth and on the moon, her/his ability to make a problem concrete is high. But if s/he doesn't measure anything, then we can show her/him concrete targets as follows:

Contact: "Let's compare values of universal gravity with those of mass on the earth and on the moon!" 32.

If s/he forgets the tool to measure mass because his/her ability to retain information in the mind is low, we can help her/him to use the tool by using the following:

Contact: "Mass of objects can be weighed by a balance beam?" 33.

After these supports, we can allow her/him to discover the difference between mass and universal gravity as:

Contact: "Were both values of mass and universal gravity fixed?" 34.

If the worker answers that the value of mass of an object was fixed but the value of universal gravity wasn't fixed, s/he has discovered the difference. But if s/he hasn't discovered it, we can show her/him a point of view as follows:

Contact: "How do the object's weight values on the moon differ between mass and universal gravity?" 35.

After these methods, we can take methods to discover rules. I would like to propose this in future work.

CONCLUDING REMARKS

I proposed educational methods for the formation and expansion of conceptions according to learners' individualities. It is my hope that, in the future, more research will be conducted to verify my hypothesis. Human individuality and learning processes are complex, therefore I only discuss a small portion of them. Future research must be conducted to further understand them.

REFERENCES

1. T.G. West, 1991. In the Mind's Eye: Visual Thinkers, Gifted People with Learning Difficulties, Computer Images, and the Ironies of Creativity, New York, Prometheus Books.
2. M. Denis, 1979. Les images mentales, Paris: Presses Universitaires de France.
3. K. Sakurai, 1987. To Creation from Discovery, Chijin sensyo24, Tokyo, Chijin Shokan. (in Japanese).
4. D.E. Segre, 1983. PERSONAGGI E SCOPERTE NELLA FISICA CLASSICA, Dalla caduta dei gravi alle onde elettromagnetiche, Milano, Arnoldo Mondadori Editore.
5. M. Monk, 1995. On the identification of principles in science that might inform research into students' beliefs about natural phenomena. International Journal of Science Education, 17(5), 565-573.
6. K. Ishino, 1989. A consideration about electromagnetism as an example of physics education. Osakayouiku University Masters Degree Thesis. (in Japanese).
7. K. Shino, K. Sugai, R. Mizoguchi, 1996. An Intelligent Education System: *Galileo* which Supports Scientific Thinking - Guidance and modeling through verbal indication of situation, Trans. Japanese Society for Information and Systems in Education, 13(1), 19-32. (in Japanese).
8. K. Shino, K. Sugai, R. Mizoguchi, 1996. An Intelligent Education System: *Galileo* which Supports Scientific Thinking - Philosophy and Basic Architecture, PRICAI'96: Topics in AI (Lecture Notes in AI 1114 and subseries of Lecture Notes in Computer Science, Heidelberg, Springer-Verlag, 71-84.
9. K. Ishino, 1997. Discussion about worker's thinking on how to move objects and support methods for it. in Proceedings IPMM'97, Gold Coast, Australia, Eds., T. Chandra, S.R. LeClair, J.A. Meech, B. Verma, M. Smith, B. Balachandran, Vol. 2, 934-940.
10. F. Leach, R. Driver, R. Millar, P. Scott, 1997. A study of progression in learning about the 'nature of science': issues of conceptualization and methodology, Inter. J. of Science Education, 19(2), 147-166.
11. K. Ishino, 1997. Consideration about expansion of conceptions in learner's mind - Formed around concepts regarding "falling objects, Japanese Society for AI - Report of Research, SIG-J-9701, 1-6.
12. A. Domenech, E. Casasus, M.T. Domenech, 1993. The classical concept of mass: theoretical difficulties and students' definitions, International J. of Science Education, 15(2), 163-173.

Robotics and Intelligent Control II

Navigation by Weighted Chance

S. Reimann*, A. Mansour**

*Institute for Autonomous intelligent Systems
GMD - German National Research Center for Information Technology
Birlinghoven, Germany
Email: stefan.reimann@gmd.de

**Bio-Mimetic Control Research Center
The Institute of Physical and Chemical Research (RIKEN), Nagoya, Japan
Email: mansour@nagoya.riken.go.jp

ABSTRACT

This paper deals with the problem of how to control the movement of a simple robot which has the goal to reach a specified target within finite time and to stay within some pre-defined distance to it. The system's design proposed is as minimal as possible and reflects the basal reflex arc as observed in biological systems. The dynamics is due to a multiplicatively modified random walk. In particular only one simple, omni-directional sensor is used so that the robot does not receive any directional information about the target. The mobile robot shows a reliable and fast homing behavior towards a defined area and stays in some given neighborhood of it. The computational effort needed is seen to be minimal.

Keywords: mobile simple robot, system control, stochastic approach, low-dimensional control, algorithm.

INTRODUCTION

The motion of simple animals, such as protozoa, bacteria, up to insects, is commonly regarded as a kind of random walk. Correspondingly, diffusion-reaction like processes have been considered in order to describe their fundamental motion patterns up to the emergence of grouping behavior (for further reading see [6]). Moreover, the assumption of random movement has made thermodynamic considerations a natural tool for analyzing systemic properties. The assumption of random movement certainly is reasonable, when considering a mobile system having a large number of degrees of freedoms, but also smaller systems with quasi-periodic or chaotic behavior [3].

We followed this idea of a random walk as a simple model for the motion of an *agent*. As a technical example, one may think about a simple **mobile robot**, which can move in a simple environment, for example an infinite, smooth plane, having a motor **E** of an appropriate number of degrees of freedom. According to classical mechanics, its spatial state is given by its spatial coordinates $\mathbf{Q} \in \mathbb{R}^d$, $d = 2, 3$ and its momentum

$\mathbf{P} \in \mathbb{R}^d$. The action of its motor is to change its spatial state due to

$$F: \begin{matrix} \mathbf{Q} \\ \mathbf{P} \end{matrix} \rightarrow \begin{matrix} \mathbf{Q}' \\ \mathbf{P}' \end{matrix} = \begin{matrix} \mathbf{Q} + \tau \mathbf{P} \\ \mathbf{D}(\alpha) \mathbf{P} \end{matrix}$$

where $\tau \in \mathbb{R}$ is some scaling constant and $\mathbf{D}(\alpha)$ denotes a rotation of the momentum around some randomly chosen angle $\alpha \in [0, 2\pi]$. The corresponding dynamics simply is a random walk in \mathbb{R}^d .

As the *agent's* sensory pole **S**, we considered an *omni-sensor*, which is sensitive to light (of some frequency). According to the intensity measured, the sensor will produce some electrical signal v , which is supposed to be transferred to the motor along the pathway $\mathbf{S} \rightarrow \mathbf{E}$. Note that accordingly the sensory signal induced does not contain any directional information.

The motor can be thought to be either autonomous or dependent: if autonomous, its activity does not depend on the signal coming from the sensor (det $D(\alpha)$ is independent of the sensory signal v). But if the motor is dependent, det $D(\alpha)$ may be some function of v . For simplicity, it is assumed that the absolute value of the robot's velocity is constant, i.e. independent of v , and the effector action exclusively consists in changing the direction of the momentum randomly.

As our key assumption, we assume that the *agent* has an internal component \mathbf{I} whose states are called the *internal states* of the *agent*. The role of internal state, *essential variables*, was already mentioned by R.W. Ashby [1,2]. A simple model for the adaptive regulation of cells by modulation of sensitivity was analyzed in [7]. More general considerations of the biological background can be found in [8].

Receptive signals are supposed to affect the *agent's* internal state according to some function g , so that for each position $Q \in \mathbb{R}^d$ corresponds an internal state $x = g(Q) \in X$. Let Q' denote the *agent's* next spatial position due to the dynamics defined above. Then the internal state x' corresponding to this new position is a function of the coordinates (Q, P) . As such the evolution of the *agent's* internal state is related to its spatial movement.

Further, we assume that there exists a set $Y \subset X$ of "essential" internal states, which is called the "homeostatic range" of the *agent*. The homeostasis condition is that during its dynamics, the internal state of the *agent* has to be kept close to this homeostatic range. As the distance measure, define the distance between the internal state $x \in X$ and the homeostatic range $Y \subset X$, i.e. $d(x, Y) = 0$ if and only if $x \in Y$, i.e. if the internal state is homeostatic. A direction P is regarded as "GOOD", if the internal state related to the new position Q' is closer to the homeostatic range than that related to the former position Q . This "weight" is formally defined as:

$$c(Q, P) = \begin{cases} +1 & \text{iff } d(x, Y) \geq d(x', Y) \\ -1 & \text{else,} \end{cases}$$

Suppose that the actual position of the *agent* is (Q, P) . Then define the forward cone as:

$$K_+ = \{p \in \mathbb{R}^d : A < \langle p, P \rangle P^2\}$$

where $\langle \cdot, \cdot \rangle$ denotes the ordinary the scalar product on \mathbb{R}^d and $0 < A < P^2$. Obviously, A is related to the opening "angle" of the cone. Analogously, K_- is defined by the property $-P^2 < \langle p, P \rangle < -A$. The modified model then is:

$$F_c: \frac{Q}{P} \rightarrow \frac{Q'}{P'} = \frac{Q + \tau P}{c(Q, P)D(\alpha)P}$$

where τ has the same meaning as before, but the rotation angle α is randomly chosen from the interval $[-\arccos(A/P^2), \arccos(A/P^2)]$.

In words: the *agent* proceeds moving in its former direction, $P' \in K_+$ if this direction is GOOD, otherwise the movement of the robot is reversed, $P' \in K_-$. The random variable α represents "internal noise" of the robot itself in that its velocity is only determined up to some extent, represented by A . If $A = P^2$, the dynamics are completely deterministic, i.e. the robot either maintains or precisely reverses its former direction.

The mapping (3) does not represent a diffusion process of Langevin type. To induce noise to a (deterministic) system, a *Langevin* force term is commonly *added*, which has to fulfill certain stochastic properties (vanishing average and δ -correlation). In contrast, the modulation of the random walk in our approach is {\em multiplicative}. By comparing the two functions (1) and (2), one immediately sees that

$$F_c(Q, P) = c_*(Q, P) \cdot F(Q, P), \quad c_*(Q, P) := 1/c(Q, P)$$

Therefore, F_c can not be written as a random walk F to which an external driving force is superimposed. In fact, the dynamics of the *agent* may be regarded as a random dynamics in a gradient field. But this gradient field is internally defined, rather than externally: By mapping $g: U \rightarrow X$, the external signaling space U is mapped to an internal space X , on which the "force term" is defined. Therefore, this force can be regarded as being "generated" by the system itself.

NUMERICAL RESULTS

In the following, only some particular properties of the dynamics of the so-defined system will be considered. The model proposed represents an extension and generalization of the approach mentioned by O.E. Holland and C.R. Melhuich [5]. A mathematical analysis of the mapping including stability of the invariant set and discussion of ergodicity are beyond the scope of this work. The main aim of this part is to very roughly compare the random walk F with the "weighted" random walk F_c .

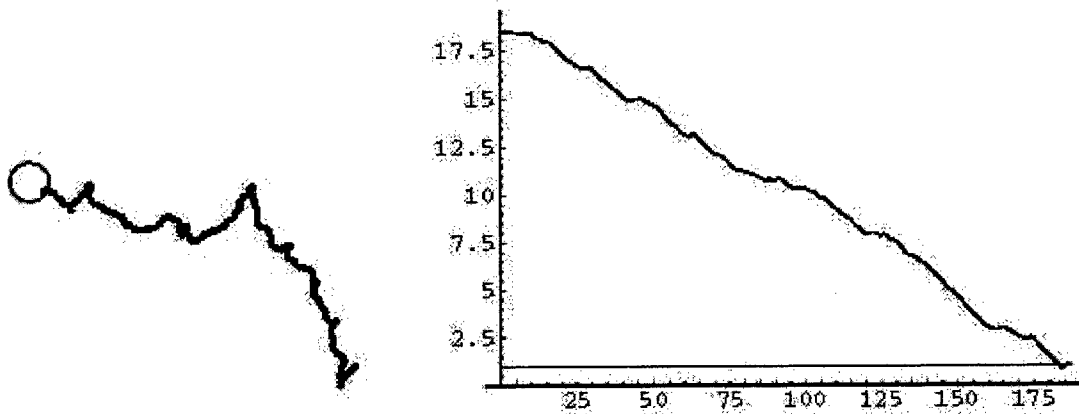


Fig. 1: Phase-plot of the weighted random motion of the agent and the time-development of its spatial distance from the source, due to Equation 3. The *agent's* initial position is far away from the source.

The above figures display the spatial motion of an *agent* due to the mapping F_c , being subject to a signaling field emitted by some source, which is located inside the circle. The simulation was done for a very simple model: as the signaling source, we defined a light bulb of constant intensity, so that the signaling field emitted is proportional to the field of light intensity. Accordingly, the strength of the receptive signals was considered as a monotonously decreasing function of the spatial position of the agent. The action of its effector, i.e. the motor, was assumed to be autonomous, leaving the velocity unchanged, $\det \mathbf{D}(\alpha) = 1$. Moreover, the action of the receptive signals on the internal state was assumed to be strictly monotonous, i.e. the mapping g was assumed to be strictly monotonous. As such, the disc displayed below reflects the homeostatic range of the agent in the spatial domain. In the simulations, the disc was assumed to have a finite extension, $0 < d < \infty$.

In contrast to a pure random walk, the trajectory due to F_c can be seen to be composed of parts of *perturbed straight lines*, so that its motion becomes directed towards the target in the mean. First of all, it is clear that no real straight lines can occur because of the random choice of the rotation matrix. But apart from this, the segments have to be lines because of our definition of the weight function, which roughly says that a certain direction is maintained under certain conditions. The appearance of "straight" trajectories also can not be expected in pure random walks.

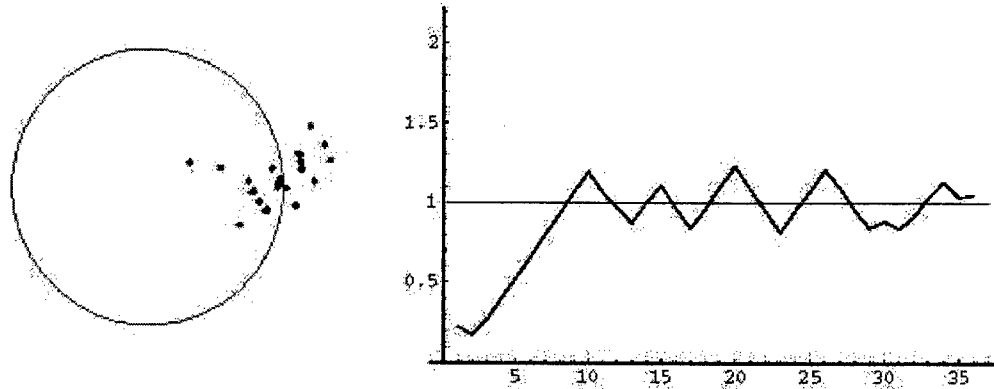


Fig. 2: Phase-plot of the weighted random motion of the agent and the time-development of its spatial distance from the source, due to equation 3. The *agent's* initial position is close to the source.

As is well-known, purely diffusive processes in a plane have a mixed property according to which each trajectory of a purely diffusive system will meet any arbitrary small neighborhood of every point in the plane after a sufficiently long time. Therefore, "homing behavior", i.e. reaching the disc, is trivially achieved by a random-walk dynamics. In fact, as apparent from Figure 2, the agent reaches the source after a couple of iteration steps, its time development being due to F_c . For the same reason, a "purely" diffusive *agent* will leave every disc of finite radius d after some time $t \approx d^{-2}$. Mapping (3) can be seen not to be mixed, in the opposite sense: According to the dynamics defined by F_c , the *agent* will not escape from a finite neighborhood of the target, but will remain in some finite distance of the disc for all time. Actually, the spatial trajectory of the *agent* shows an oscillation around the border, its amplitude being dependent on the initial velocity of the agent and the time scale parameter. This, in fact, constitutes a major difference between the pure random walk F and our "weighted random walk" F_c .

CONCLUSION

The model proposed represents a multiplicative modulation of a simple random walk: At each time-step, the direction is chosen randomly from a forward- or a backward cone according to the actual direction, due to the value of a weight function.

The forcing term is internal, rather than external, i.e. no external force field is superimposed. This "weighted" random walk exhibits dynamical aspects, which fundamentally differ from those of normal random walks:

1. The dynamics of the model proposed is not "mixed", but rather establishes a motion which converges to some given target.
2. The mean direction is directed towards the target.

The computational effort needed is seen to be minimal and, in particular, does not include any "orientation", i.e. no direct directional information is present. This process of "weighted diffusion" may be important as a framework for describing and analyzing the motion pattern of simple animals.

ACKNOWLEDGMENT

This paper represents work done in a joint project between the GMD, Germany, and the BMC - RIKEN, Japan, on the stochastic control of minimal robots. The authors would like to express their gratitude to Professor Dr. N.Ohnishi (RIKEN) and Professor Dr. Th. Christaller (GMD) for support.

REFERENCES

1. R.W. Ashby, 1952. Design for a brain, John Wiley and Sons.
2. R.W. Ashby, An Introduction to Cybernetics. John Wiley, New York;
3. C. Beck, F. Schlogl, 1993. Thermodynamics of chaotic systems, Cambridge Nonlinear Science Series 4, Cambridge University Press.
4. C.W. Gardiner, 1993. Handbook of Stochastic Methods for Physics, Chemistry, and the Natural Sciences, Springer, New York.
5. O.E. Holland, C.R. Melhuish, 1996. Getting the Most from the Least: Lessons for the Nanoscale from Minimal Mobile Agents, ALIFE V.
6. J.D.Murray, 1989, Mathematical Biology, Springer Verlag.
7. S. Reimann, 1998. Oscillation and pattern formation in a system of self-regulating cells, Physica D, 114, 338-361.
8. S. Reimann, 1998. Stability and equilibrated structure, Proceeding of the 3rd German Workshop on Artificial Life, (ed. C. Wilke, S. Altmeyer, Th. Martinez), Harry Deutsch Verlag.

Vehicle Routing Problem Using Clustering Algorithm by Maximum Neural Networks

N. Yoshiike*, Y. Takefuji*

Graduate School of Media and Governance,
University of Keio, Kanagawa, Japan

ABSTRACT

The vehicle routing problem (VRP) is one of the well known optimization problems. It is used to minimize the total length of all routes of vehicles where each of the vehicle has a capacity constraint respectively. This paper proposes a self-organization neural network model for obtaining the best solution for VRP. Our method consists of two phases. In the first phase, the customers are grouped to several delivery areas for vehicles assignment by Maximum Neuron model. In the second phase, the TSP in each area is solved by Elastic net model proposed by Andrew et. al. The clustering algorithm used in the first phase is a Maximum Neuron model. Maximum Neuron model is one of the neural networks proposed by Hopfield that can minimize a cost function considering various constraints. In the second phase, Elastic net model is used to solve the problem and it can obtain good solutions of TSP. Our method improves the precision of solution, and can be extended for big size problem. Our simulation result shows that Maximum Neuron model can achieve better solutions than other methods in certain conditions.

INTRODUCTION

The vehicle routing problem[3] (VRP) can be stated as follows: A set of L vehicles, with same capacity Q , is located at depot D . City i is located at X_i in the two dimensional map and has a demand q_i . Each vehicle, finding a route which begins at the depot, visits a subset of cities and returns to the depot without violating the capacity constraint. The objective is to minimize the total length of all routes. (see Fig. 1.).

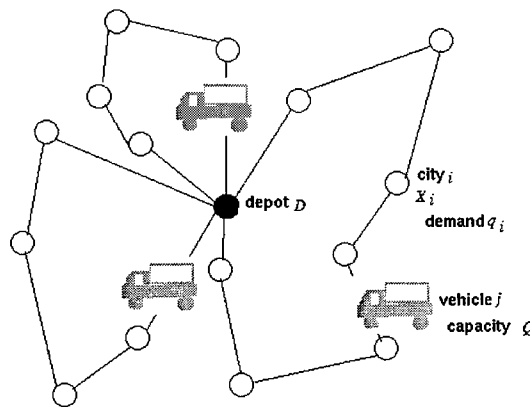


Fig. 1: Vehicle Routing Problem

This problem is widely applied to many real life delivery problems. Delivery operations of goods to and from customers is one of situations where this problem can be applied to. The collection of mail from mail-boxes and the operation of school bus services are well known examples of deliveries from customers.

The Elastic net algorithm for solving the VRP is proposed in [1]. L sets of dynamic *rubber band*, which is initialized as L small loops, is stretched towards the cities with elastic forces of the band. Although the

Elastic Net algorithm will find good solution for the traveling salesman problem (TSP), it is difficult to obtain the best solution for the VRP in realistic computational time.

This paper proposes a clustering algorithm using the Maximum Neuron model to simplify the VRP to the TSP. By defining the independent routing cost for each customer, VRP can be divided into two problems: the vehicle assignment problem and the TSP. In the first phase, the vehicle assignment problem is solved by the Maximum Neuron model, and in the second phase, the TSP is solved by the Elastic Net model.

CLUSTERING ALGORITHM

The cost c_{ij} for assigning vehicle j to city i can be defined as distance from *main route* of vehicle j to city i .

c_{ij} can be described as follows:

$$c_{ij} = \begin{cases} |N_{ij} - X_i| & \text{if } (M_j - D) \cdot (X_i - D) > 0 \\ |M_j - X_i| & \text{otherwise} \end{cases} \quad 1.$$

M_j is initially set to the centroid of the cities visited by vehicle j and updated by iteration and N_{ij} is the position of the foot of the perpendicular from X_i of the line L_j which passes through D and M_j . The cost c_{ij} is the distance from L_j to X_i if X_i locates the side of M_j or the distance from M_j to X_i otherwise (see Fig. 2).

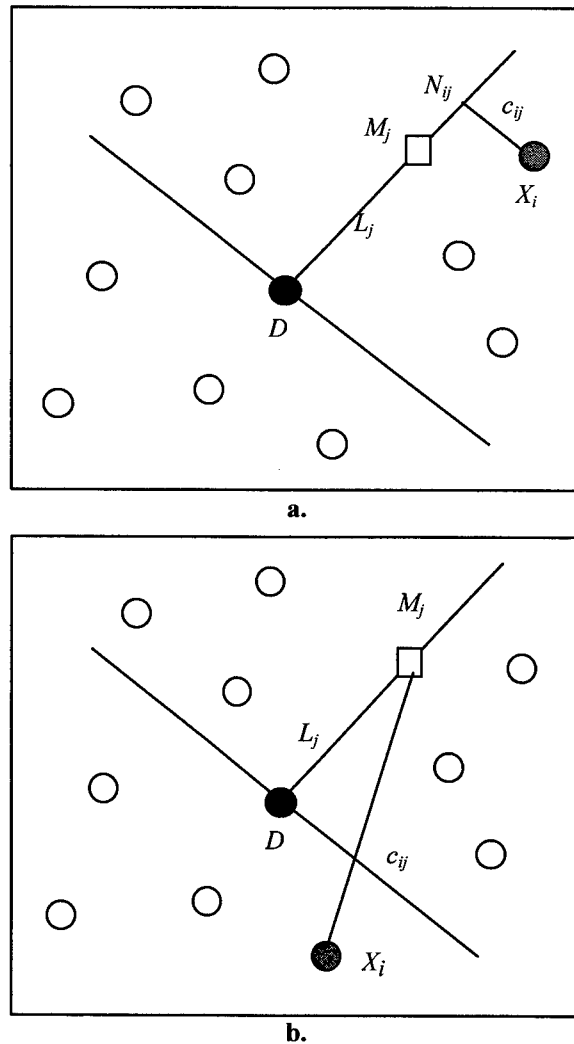


Fig. 2. A cost for visiting city i by vehicle j . a. when X_i locates the side of M_j . b. otherwise.

V_{ij} indicates whether vehicle j visits city i . $V_{ij} = 1$ if vehicle j visits city i and $V_{ij} = 0$ otherwise. V_{ij} must satisfy Eq. 2. and Eq. 3.

$$\sum_j V_{ij} = 1 \quad 2$$

$$\sum_i V_{ij} \cdot q_i < Q \quad (\forall_j) \quad 3.$$

The total cost can be described as follows:

$$\sum_j \sum_i c_{ij} \cdot V_{ij} \quad 4.$$

NEURAL NETWORK DYNAMICS

An approximate optimal solution of V is searched for by using the Maximum Neuron model [8]. To satisfy constraint 2 constantly, V_{ij} is calculated by:

$$V_{ij} = \begin{cases} 1 & \text{if } U_{ij} = \max \{U_{ia}\} \quad (a = 1, \dots, L) \\ 0 & \text{otherwise} \end{cases} \quad 5.$$

where U , initially set by random number selection, is updated by:

$$U_{ij}(t+1) = U_{ij}(t) + \frac{dU_{ij}}{dt} \quad 6.$$

$$\frac{dU_{ij}}{dt} = -\alpha \cdot c_{ij} - \beta \left(\sum_a V_{aj} \cdot q_a - Q \right) \quad 7.$$

α is the coefficient of cost c_{ij} , and β is the coefficient of the capacity constraint of a vehicle. α approaches 0 from an initial value by the following equation:

$$\alpha(t+1) = \alpha(t) - \frac{\alpha(0)}{T} \quad 8.$$

where $\alpha(0)$ is the initial value of α .

M_j , which approximates the centroid of cities visited by vehicle j , is initialized by:

$$M_j(0) = \frac{\sum_a X_a V_{aj}}{\sum_a V_{aj}} \quad 9.$$

where $M_j(0)$ is the initial value of M_j and is updated by:

$$M_j(t+1) = M_j(t) + \gamma \left\{ \frac{\sum_a X_a V_{aj}}{\sum_a V_{aj}} - M_j(t) \right\} \quad 10.$$

γ is a learning parameter of M_j , and approaches 0 from an initial value based on the following equation:

$$\gamma(t+1) = \gamma(t) - \frac{\gamma(0)}{T} \quad 11.$$

where T is the frequency for updating α and γ . V_{ij} is updated by Eq. 5. until V satisfies constraint 3 and until the change in total cost equals zero.

SIMULATION

In the second phase, the problem(TSP) is solved by an Elastic Net model[2, 4]. We approach eil50, eilA76, eilA101 and problem4 in the TSPLIB[9] for our simulation. The simulation results of these problems are shown in Fig. 3., where the route of each vehicle is drawn by different line colors.

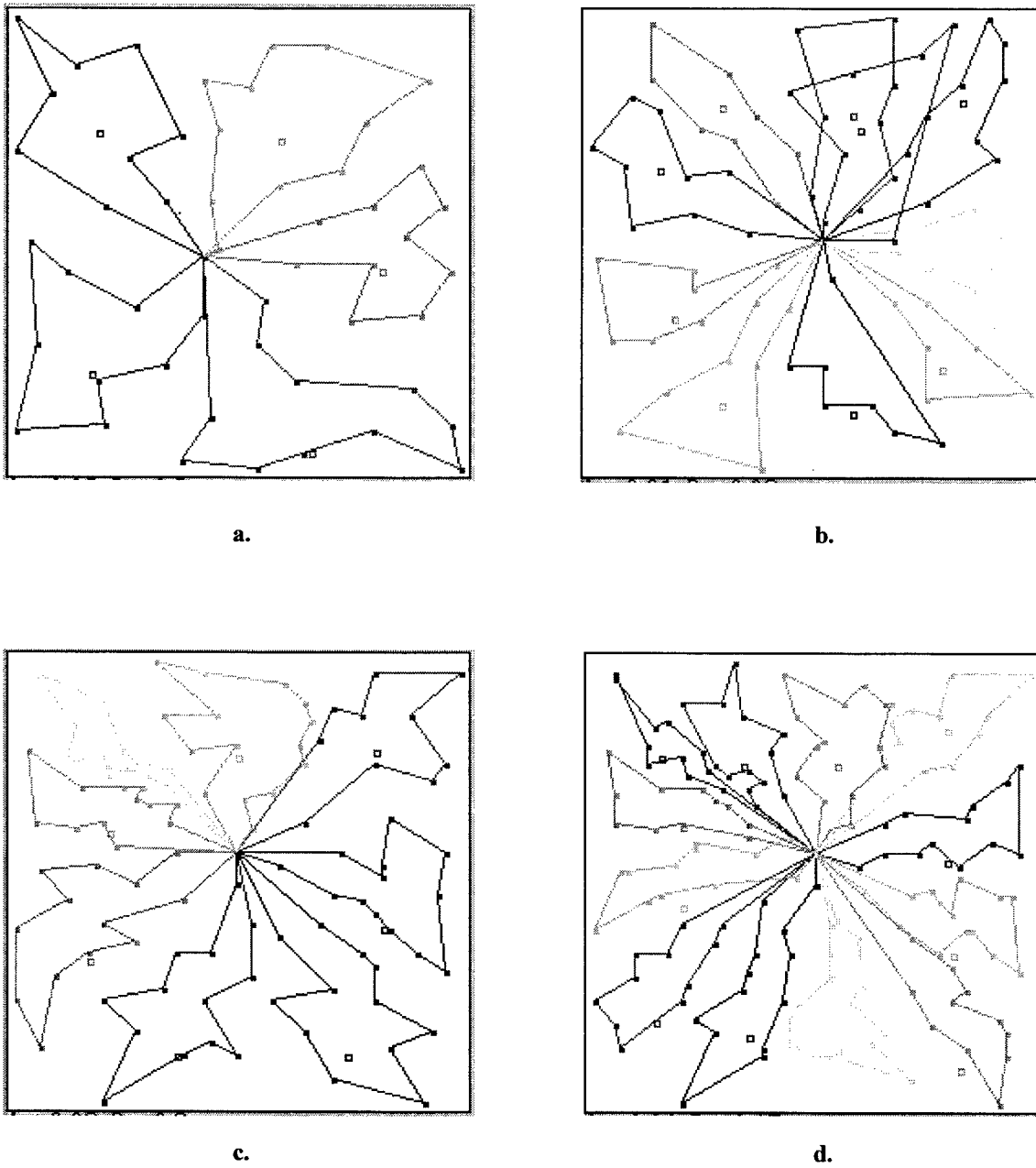


Fig. 3: Simulation results of eil51, eilA76, eilA101 and problem4.

- a.** $\alpha = 0.02, \beta = 0.30, \gamma = 0.03$, step = 299, number of vehicles = 5, total distance = 522
- b.** $\alpha = 0.01, \beta = 0.08, \gamma = 0.03$, step = 302, number of vehicles = 10, total distance = 898
- c.** $\alpha = 0.02, \beta = 0.30, \gamma = 0.03$, step = 300, number of vehicles = 8, total distance = 862
- d.** $\alpha = 0.04, \beta = 0.15, \gamma = 0.03$, step = 300, number of vehicles = 12, total distance = 1099

The Maximum Neuron model solutions are compared with the best known solutions in Table 1.

Table 1: Results of the Maximum Neuron on the VRP.

problem	number of cities	total distance					number of vehicles				
		A	B	C	D	E	A	B	C	D	E
eil51	50	521	560	585	556	534	5	6	6	5	-
eilA76	75	898	-	900	876	871	10	-	10	10	-
eilA101	100	853	-	887	863	851	8	-	8	8	-
problem4	150	1081	-	1204	-	1064	12	-	-	-	-

A = Maximum Neuron model; B = Elastic Net [1]; C = The savings approach [5]; D = 3-optimal method [6]; E = Tree Search [7].

DISCUSSION

The experimental results show that our method can obtain better solution than all other methods, B, C, D and E, in eil51. Comparing with the Elastic Net model, our solution for eil51 gives less vehicles than that of the Elastic Net model, and the total distance is approximately 7% better. By separating VRP into two problems, the method improves the accuracy of solution, and it can also perform when the number of cities increases. However, with 75 customers, our simulation result does not show the best solution. Because this method is heuristic in nature, approximating the cost function in the first phase, it is not guaranteed to find the best solution for all problems. But the experimental results show that the iterations do not increase when the number of cities increases. Since U_{ij} can be calculated in parallel by the algorithm described by Eq. 7., the model is able to work in parallel processors. Therefore, if our model is simulated by parallel processors, it can find solutions in $O(n)$ of computational time where n is the number of the cities.

CONCLUSION

By deciding on an independent routing cost for each customer, we can divide VRP into two problems. Using a clustering method in the first phase, VRP can be simplified to a kind of TSP. Our method improves the precision of the solution, and can be extended to larger size problems. Our method can obtain approximately optimal solutions in the fixed iterational steps. Our method is particularly good at obtaining the best solution for eil51 which has 50 cities.

REFERENCES

1. A.I. Vakhutinsky, B.L. Golden, 1994. Solving Vehicle Routing Problems Using Elastic Nets, IEEE Inter. Conf. on Neural Networks. IEEE World Cong. on Computational Intelligence, 7, 4535-4540.
2. R. Durbin, D. Wilshaw, 1987. An Analogue Approach to the Traveling Salesman Problem using an Elastic Net Method, Nature 326, 689-691.
3. S. Eilon, C. D. T. Watson-Grandy, N. Christofides, 1971. Distribution Management, Hafner Publishing Company, New York.
4. M. W. Simmen, 1991. Parameter Sensitivity of the Elastic Net Approach to the Traveling Salesman Problem, Neural Computation 3, 363-374.
5. G. Clarke, J. Wright, 1964. Scheduling of Vehicles from a Central Depot to a Number of Delivery Points, Oper. Res., 12(4), 568-581.
6. N. Christofides, S. Eilon, 1969. An algorithm for the vehicle dispatching problem, Opl. Res. Q., 20, 309.
7. N. Christofides, A. Mingozzi, P. Toth, 1978. The Vehicle Routing Problem, Urbino Working Paper, July.
8. Y. Takefuji, K.-C. Lee, H. Aiso, 1992. An artificial maximum neural network: a winner-take-all neuron model forcing the state of the system in a solution domain, Biological Cybernetics, 67, 243-251.
9. <http://www.iwr.uni-heidelberg.de/iwr/comopt/soft/TSPLIB95/TSPLIB.html>

Acquisition of Communication Protocol for Autonomous Multi-AGVs Driving

M. Watanabe*, M. Furukawa, Y. Kakazu*****

*Information Processing Center, Asahikawa National College of Technology,
#1-6 Shunkodai 2-2 Asahikawa Hokkaido 071-8142, Japan

**Dept. of Information Technology Integration, Asahikawa National College of
Technology, #1-6 Shunkodai 2-2 Asahikawa Hokkaido 071-8142, Japan

*** Faculty of Engineering, University of Hokkaido,
Nishi 8 Kita 13 Kitaku Sapporo Hokkaido 060, Japan

ABSTRACT

In order to realize autonomous AGV driving, two problems arise. The first problem is known as the autonomous vehicle navigation problem. The second is the problem of acquiring a communication protocol. The communication protocol means that an AGV can understand other AGV intention by its behavior. If the AGV can understand the behavior of another AGV when they meet, it becomes possible for them to autonomously avoid collision and collaborate on their tasks.

This paper proposes an acquisition method for communication protocol using Q-learning. Multi-AGVs are defined by a set of agents and each agent is defined by a learning automaton with a Q-learning function. Through many experiences in transporting workpieces within a virtual factory, the AGVs gradually acquire a common behavior by Q-learning and finally adopt a suitable communication protocol. Numerical experiments are executed to examine what kind of communication protocol is acquired by Q-learning by the multi-AGVs operating together. The acquired communication protocol selects whether the AGV will "turn right" when two AGVs drive in opposite directions and meet each other. When the AGVs drive in the same direction, the acquired communication protocol selects "go forward" by maintaining the same speed for each AGV. By comparing the acquired communication protocol with other options, the AGVs that adopt the correct communication protocol, process the given workpieces slightly faster. This result shows that an AGV can acquire a suitable communication protocol by learning in any environment.

INTRODUCTION

Several concepts for a new production system have been proposed for a past few years. Typical concepts are an agile production system, a lean production system, a biological production system and a network production system. These systems direct to an autonomous decentralized production system. An intelligent material handling system makes a great role of such systems. An AGV (an automatically guided vehicle) occupies a center of the intelligent material handling system and it shoulders a main schedule part to the production system. It can flexibly carry workpieces between an automated warehouse, machine tools, assembly machines, and testing machines.

In order to drive AGVs autonomously, two problems arise. The first problem is known as the autonomous vehicle navigation problem, which has been studied for a long time.

The second problem is acquisition of a communication protocol. This treats the problem of how the AGVs can understand other AGVs intention by their behavior. If the AGV can understand other AGVs behavior when they meet on their ways, it becomes possible for the AGVs to autonomously avoid their collision between themselves and collaborate with their behavior. We call this understanding the communication protocol.

The "flowers and bees" problem [1] is known as a communication protocol acquisition problem which has been well-studied. In this problem, mutual understanding between bee actions and flower signals are taken into account. When a bee selects a suitable flower signal, it will move toward the flower and collect nectar. Also, when a flower emits an attractive signal towards the bees, it can guide a bee to itself and have the bee carry its pollen. This mutual understanding can be considered as a communication protocol. One difference between the AGV and "flowers and bees" problems is that the latter problem deals with heterogeneous agents while the AGV communication problem considers homogeneous ones. As well, the flowers and bees problem is a static problem while AGV driving is a dynamic one.

This paper proposes a communication protocol acquisition method which uses Q-learning. Multi-AGVs are defined by a set of agents and each agent is defined by a learning automaton with a Q-learning function. Multi-AGVs behave like a mob. Directions to move for each AGV are "go forward", "turn right" and "turn left". Through many experiences in transporting workpieces in a virtual factory, the AGVs gradually generate a common behavior by Q-learning and finally acquire a suitable communication protocol. This process seems to be like mob-learning.

Numerical experiments have been executed to examine the kind of communication protocol acquired by Q-learning through multi-AGVs driving. The acquired communication protocol selects an AGVs' behavior to move "turn right" when two AGVs drive in opposite directions and meet each other. On the other hand, when AGVs drive in the same direction, the acquired communication protocol selects "go forward" by maintaining the same speed for each vehicle. By using the acquired protocol, an AGV driving simulation is executed for transporting given workpieces in a virtual factory. The results show that the AGVs which adopt the acquired communication protocol process the given workpieces slightly faster than those AGVs which adopt other methods of communication.

MULTI-AGVs MODEL DEFINED BY LEARNING AUTOMATA

Automata are applied to model multi-AGVs. The multi-AGVs are defined by a set of automaton as follows:

$$A = \{A_i, i = 1, 2, \dots, n\} \quad 1.$$

where n is a total number of the AGVs. Each AGV is defined by an automaton as follows:

$$A_i = (I_i, O_i, S_i, F_i, G_i) \quad 2.$$

where I_i , O_i , S_i , F_i , and G_i are input, output, state, state action function, and output function, respectively. The state action function and output functions are represented as a function of Q-value that is defined in term of Q-learning. In general, knowledge of the AGV consists of two parts as shown in Fig. 1. The first part represents a clear knowledge that can be set up in advance with a definite knowledge. The second part represents an acquired knowledge in adaptation and learning. The model defined by Eq. 2 corresponds to the second part.

Constraints

A minimum number of sensor and communication methods are assumed to be available to each AGV. The following constraints are introduced.

- (1) An AGV can sense another AGV moving in a particular direction from its present position.
- (2) An AGV can sense another AGV moving in a direction adjacent to itself.
- (3) The machining time for a conveyed workpiece is 0 and the buffer size of an AGV being stored at the machine tool location is infinity.

These constraints were introduced based on Kaibara analysis to compare our experimental results with his.

Input

Input I_i is defined as follows:

$$I_i(t) = \{P_i(t), M_i(t), P_j(t), D_j(t)\}$$

3.

where t expresses time. Each entity of the input equation is defined as follows:

$P_i(t)$ is the current position, (x_i, y_i) of the AGV A_i .

$M_i(t)$ is the machines tool location, (x_m, y_m) to which that the current AGV is conveying the workpieces.

$P_j(t)$ is another AGV's position, (x_j, y_j) sensed by the AGV A_i in its neighborhood.

$D_j = \{N, E, S, W\}$ is the direction of movement of the sensed AGV A_j , where N, E, S, and W indicate north east, south, and west respectively.

Output

Output O_i is an action taken by AGV A_i . It is defined as:

$$O_i = \{F, R, L\}$$

4.

where F, R, and L express "go forward", "turn right" and "turn left" respectively for the action of the AGV at time $(t-1)$.

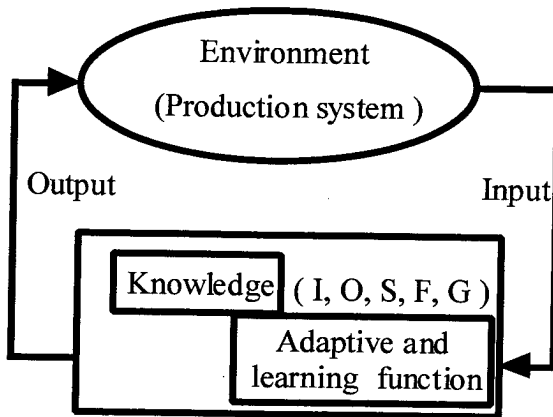


Fig. 1. learning automaton model for an AGV.

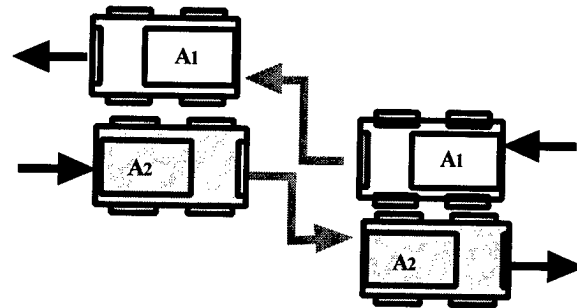


Fig. 2. Collision Avoidance when two AGVs meet.

STATE

State $S_i(t)$ of the AGV A_i at an instant t is composed of the current moving direction $D_i(t)$, its previous moving direction $D_i(t-1)$ and the action $O_j(t)$ of the sensed AGV A_j . Namely, it is described by:

$$S_i = \{D_i(t), D_i(t-1), O_j(t); i \neq j\}$$

5.

STATE ACTION FUNCTION AND OUTPUT FUNCTION

An automaton usually uses a state transition function. In order to incorporate a learning function into the automaton, a state action function is adopted instead of a state transition function. The state action function F_i maps a set of states into a set of outputs. The mapping is performed through a Q-value in terms of Q-learning. An output function G_i usually maps a set of states into a set of actions. However, the output function G_i is changed into one that maps one set of actions into another set of actions because of adopting Q-learning. This output function is called as the action selector. These two functions will be defined later.

ACQUISITION OF A COMMUNICATION PROTOCOL BY LEARNING

If an AGV can understand the behavior of another AGV when they meet, this understanding makes it possible to avoid collision and cooperate in their work. This mutual understanding is called n-formulated in this section.

DESCRIPTION OF THE PROBLEM

Let a set of two AGVs be defined as A_i and A_j , respectively. An action of A_j makes a state of A_i , while an action of A_i makes a state of A_j with the model defined above. Let set C_i and C_j be the profit which depends on actions of A_i and A_j , respectively. Then, the acquisition problem of a communication protocol is to determine the mutual actions of A_i and A_j so as to maximize the total profit of C_i and C_j when they meet. This problem is formulated as follows.

$$\begin{aligned} & \max_{O_i, O_j} (C_i + C_j) \\ & \text{subject to} \\ & O_i = G(O_i(S_i)) \text{ and } O_j = G(O_j(S_j)) \end{aligned} \quad 6.$$

Here, C_i and C_j can be set as 0 when each action makes a collision of A_i and A_j and as some positive values when each action avoids a collision of them. Since the states S_i and S_j contain the actions of A_j and A_i , respectively, the above formulation is rewritten as:

$$\begin{aligned} & \max_{O_i, O_j} (C_i + C_j) \\ & \text{subject to} \\ & O_i = G(O_j) \text{ and } O_j = G(O_i) \end{aligned} \quad 7.$$

The actions O_i and O_j are determined through a Q-value acquired by Q-learning. Then, this problem can be treated as a learning problem through experience.

ACQUISITION OF A COMMUNICATION PROTOCOL BY Q-LEARNING

Q-learning is a learning method, which determines suitable actions by estimating a future profit by trial and error. Let us assume that two AGVs A_i and A_j are adjacent to each other. If actions O_i and O_j are selected for A_i and A_j respectively, then the moving directions $D_i(t)$ and $D_j(t)$ of A_i and A_j are determined respectively. By using current positions $P_i(t)$ and $P_j(t)$ of A_i and A_j , the next position $P_i(t+1)$ and $P_j(t+1)$ can be calculated by $P_i(t+1) = P_i(t) + D_i(t)$ and $P_j(t+1) = P_j(t) + D_j(t)$, respectively. Based on the predicted positions $P_i(t+1)$ and $P_j(t+1)$, profits C_i and C_j are estimated by the following algorithm in the cases where two AGVs drive in the same direction ($D_i(t) = D_j(t)$) and in the opposite direction ($D_i(t) \neq D_j(t)$), respectively:

```

If  $P_i(t+1) \neq P_j(t+1)$  and  $P_i(t+1) \neq \forall P_k(t)$ 
    then  $C_i = \text{Reward}$  and  $C_j = \text{Reward}$ 
else if  $P_i(t+1) \neq P_j(t+1)$  and  $P_i(t+1) \neq \forall P_k(t)$ 
    then  $C_i = \text{Penalty}$  and  $C_j = \text{Penalty}$ 
else
     $C_i = \text{Penalty}$  and  $C_j = \text{Penalty}$ 
endif

```

In the above algorithm, Reward and Penalty are positive and negative values respectively, and the subscript k means an arbitrary AGV around A_i and A_j . In the above algorithm, the first judgement corresponds to the success of collision avoidance between A_i , A_j and an adjacent AGV A_k . The second judgement corresponds to the failure of collision avoidance. The Q-value, when the state of A_i is S_i and its action is O_i , is updated by use of the profit C_i as a direct cost of Q-learning as follows:

$$Q(S_i, O_i, t+1) = (1-\alpha)Q(S_i, O_i, t) + \alpha[C_i + \gamma \max_h Q(S_i, O_h, t+1)] \quad 8.$$

where α and γ are learning parameters. This equation becomes the state action function F_i .

The action O_i of A_i is determined by use of a Boltzman-distribution function as follows:

$$O_i = \underset{h}{\text{argProb}}_{O_h} \left[\frac{\exp\{Q(S_i, O_h, t)/T\}}{\sum_h \exp\{Q(S_i, O_h, t)/T\}} \right] \quad 9.$$

where **Prob** is a function that selects action O_h depending on its probability, **arg** is a function to extract the subscript of O_h , and T is a computational "temperature".

EXPERIMENT FOR ACQUIRING A COMMUNICATION PROTOCOL

Simulations were performed to acquire a communication protocol by Q-learning as formulated above. Also, simulations were done to verify that acquired communication protocol works well.

ACQUISITION OF A COMMUNICATION PROTOCOL

Figures 3 and 4 show graphical output of the numerical experiments. A virtual factory is considered in which three machine tools, M_1 , M_2 and M_3 , are located on a horizontal centerline. 100x100 cells represent the field of the factory. The size of an AGV is represented by one cell. A hundred AGVs are generated at random positions within the factory and they convey workpieces to the machine tools in accordance with the predetermined order. The machining order is M_1 , M_2 and M_3 or M_3 , M_2 and M_1 . One order is randomly given to an AGV. It is assumed that the machining time for the given workpieces is 0 and the buffer size to store AGVs at a machine tool is infinite. These conditions are the same as Kaibara's experimental ones.

Parameters α and γ for Q-learning are set to 0.3 and 0.8, respectively. Also, values of Reward and penalty are set to 0.1 and -0.01, respectively. It is counted as one cycle that a hundred of AGVs have finished to convey workpieces to the machine tools in accordance with the specified order. Learning is repeated for up to forty thousand cycles and the Q-value is updated on each cycle.

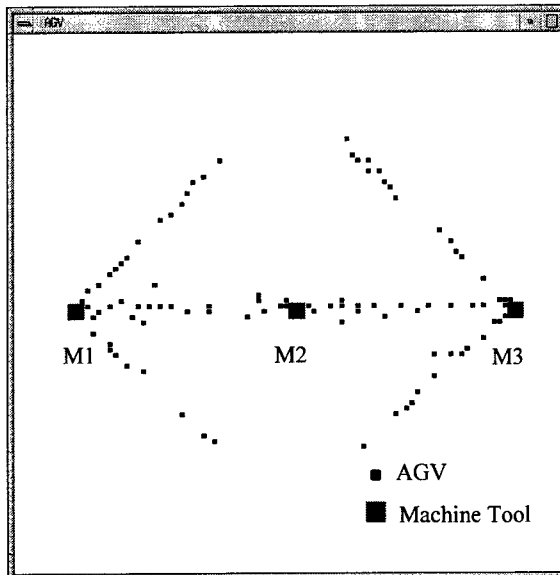


Fig. 3. AGVs driving simulation for learning at an early step.

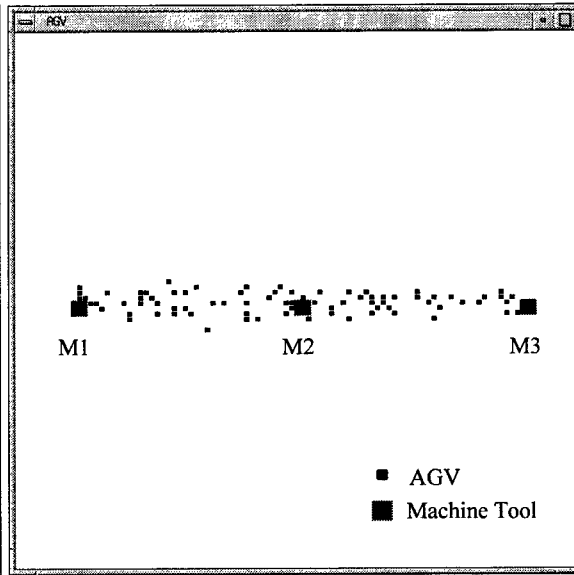


Fig. 4. AGVs driving simulation for learning at an intermediate step.

Fig. 5 shows a learning curve for two AGVs drive in the same direction while Fig. 6 depicts the case where two AGVs drive in the opposite direction. It is seen that the Q-value corresponding to "go forward" is acquired in a monotonously increasing fashion when the AGVs drive in the same direction. On the other hand, the Q-value corresponding to "turn right" is acquired when the AGVs drive in opposite directions. In the latter case, other actions are acquired to a small extent. When the AGVs drive in opposite directions, it is observed that the Q-value corresponding to "turn left" is mostly acquired and, other actions are acquired at a small extent in other experiments. From the results, the desired protocol is to take an action "go forward" when the AGVs drive in the same direction and to mostly take an action "turn right" or "turn left" when the AGVs drive in the opposite direction.

CONVEYING WORKPIECES

The acquired communication protocol is compared to the one proposed by Kaibara [2] and to a heuristic protocol which uses a job completion ratio (OCR). An OCR means the average number of machine tools to which all AGVs convey workpieces. The heuristic communication protocol is that the AGV selects the action "go forward" when the AGVs drive in the same direction and to select the action "turn right" when the AGV drive in the opposite direction. Kaibara's protocol uses the following procedure:

```

If  $D_j(t) = D_j(t+1)$ 
  then move  $A_i$  and  $A_j$  forward respectively
if  $D_j(t) \neq D_j(t+1)$ 
  then
    If  $P_i(t+1) \neq P_j(t+1)$  and  $P_i(t+1) \neq \forall P_k(t)$ 
      then move  $A_i$  and  $A_j$  forward respectively
    else if  $P_i(t+1) = P_j(t+1)$  and  $P_i(t+1) \neq \forall P_k(t)$ 
      then move  $A_i$  and  $A_j$  at random
    else
       $D_i(t) = D_j(t+1)$  and  $D_j(t+1) = D_i(t)$ 
  endif
endif
endif

```

Two AGVs are driving in the same direction when the first judgement is satisfied. On the other hand, two AGVs driving in the opposite direction satisfies the second judgement. Figure 7 shows a comparison of the OCR between the three protocols. The same model as shown in Figure 3 is used for the experiment except that five machine tools are located on the centerline. Figure 8 shows the result when ten machine tools are placed in the factory. The OCR curve of the acquired communication protocol becomes very similar to that of the heuristic communication protocol. Kaibara's protocol produces the worst result. In Figures 7 and 8, a step means an elapsed time. The AGVs using the acquired protocol complete the conveying job 20 steps faster and 8 steps faster than those using the heuristic protocol in Figures 7 and 8, respectively.

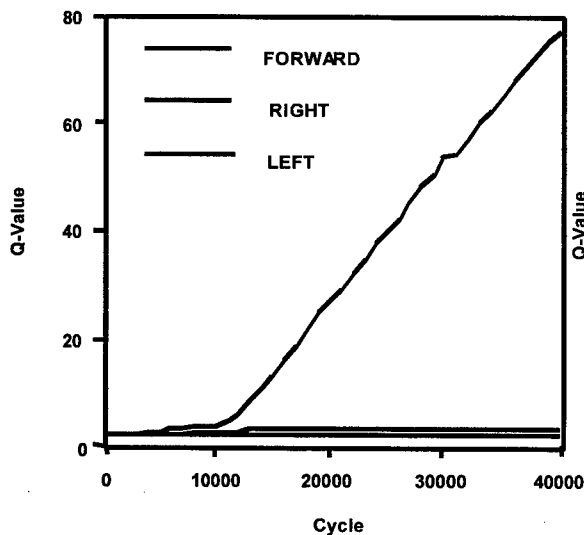


Fig. 5. An acquired protocol when two AGVs drive in the same direction.

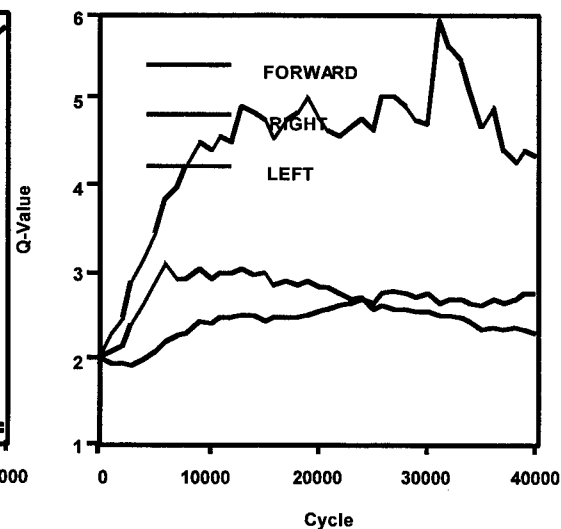


Fig. 6. An acquired protocol when two AGVs drive in the opposite direction.

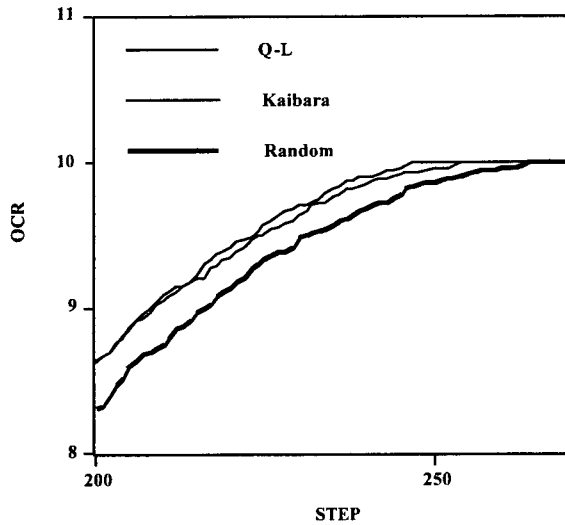


Fig. 7. A comparison of OCR between Q-learning, Kaibara's protocol and a heuristic one. Five machine tools in the factory.

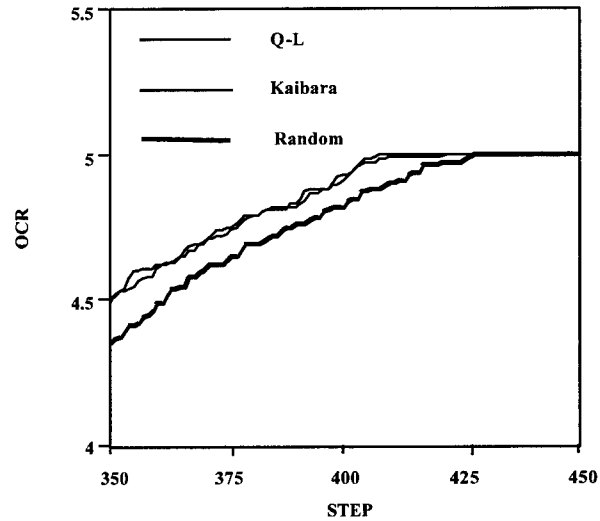


Fig. 8. A comparison of OCR between Q-learning, Kaibara's protocol and a heuristic one. Ten machine tools in the factory.

CONCLUDING REMARKS

An acquisition method for communication protocols based on Q-learning is proposed and the acquired protocol is examined by comparing it to other protocols. The following conclusions can be made:

- (1) The acquisition of a communication protocol for autonomous multi-AGVs has been formulated as a learning problem and the use of Q-learning is proposed.
- (2) The acquired communication protocol well compares to the heuristic protocol, which is intuitively given by a human. The two methods are better than Kaibara's protocol.
- (3) The AGVs that adopt the Q-learning method are freed from the dead-lock phenomenon because of the probabilistic action selection.

ACKNOWLEDGEMENT

This research is partially supported by the Japanese Ministry of Education, Science, Sports and Culture under the Grant-in- Aid for the Scientific Research (Basic research C, Account No. 412).

REFERENCES

1. N. Ono, T. Ohira, A.T. Rahm, 1998. Emergent organization of interspecies communication in Q-learning. Proc. of European Conference on Artificial Life, 396-405
2. T. Kaihara, S. Fujii, S. Kunimasa, 1998. An evolutionary self-organization scheduling using coordinated autonomous work agents. Proc. Japan-U.S.A. Sym. on Flexible Auto., III, 1375-1381

Heuristic Neuro-Fuzzy Model for Evaluation of Urban Transportation Projects

Marcus V. Quintella Cury, Saul Fuks

Coordenação de Pesquisa e Pós-Graduação em Engenharia - COPPE,
Universidade Federal do Rio de Janeiro - UFRJ, Rio de Janeiro, Brazil

ABSTRACT

Urban transportation planning projects are largely concerned with providing two benefits: trip time reduction and operating cost reduction; both of which are usually quantified in monetary form. However, neglect of *human* factors in transportation projects and failure to recognise urban transportation as a basic need of a society, have resulted in systems that do not execute their primary objective: to move human beings quickly, comfortably, economically and effectively. As well, a poorly implemented transportation system can impact significantly on the environment and affect the psychological and physical health of a population, whether they be system users or not. Many existing transportation systems fall into the above classes because they were not designed initially from the viewpoint of improving the quality of life.

In most cases, quantitative methods are preferred methods for system evaluation, with qualitative concerns and intangible benefits, perhaps acknowledged, but rarely accounted for during decision-making. Single-valued criterion play a major role in economic models in which all input factors are expressed monetarily or by applying weights of importance. There is an obsession to quantify variables in monetary terms. Exact quantification is usually impossible or unreliable since many measures must be assumed or guessed.

As an alternative, the work described here develops a methodology to evaluate urban transportation systems based on qualitative descriptions, wherein humans and their environment are the main factors. The technique tries to analyse elements such as comfort, safety, speed of movement, satisfaction, convenience, environmental pollution, social effects, economic synergism, physical and psychological effects and so on.

The model contains a knowledge base, composed of physical, chemical, biological, technical, economic and political data used to provide a humanistic perspective. The method is based on a neuro-fuzzy approach with heuristic rules organised into a hierarchy designed into a participatory environment of experts, current users and future users. All inputs are transformed into linguistic variables by ascribing degrees of belief to arrive at a final conclusion to assign a Humanistic Grade (HG) to the overall analysis. So all qualitative and quantitative inputs are compressed into a single parameter in this approach represented as HG. Since there are few mathematical methods for this approach, neuro-fuzzy technology has been chosen because of its ability to mimic human decision-making and because of the ease of model-validation.

The architecture proposed is composed of several fuzzy associative maps or sub-systems, wherein input variables are fuzzified to establish degrees of belief; the fuzzy linguistic terms are inserted into fuzzy rules which infer intermediate linguistic sub-goals. These internal variables unite with other sub-goals that feed into additional fuzzy rulebases. This process continues until a final linguistic variable describing HG is found. This variable is passed through a defuzzification process to determine a crisp output for HG, ranging from 0 to 1. In this way, HG can be used to rank projects or decide on acceptance or rejection of a project.

The knowledge contained within fuzzy sets and rules has been taken from interviews with experts. Much of this information comes from human experience, and often, is not supported by theory or empirical data. This wisdom is basically informal judgements of domain experts known as heuristic knowledge. Heuristic-thinking involves searching the problem space, learning about facts in the domain, being able to explain information and decisions, and repetition of this process until the problem is resolved. So the proposed model represents a *collective mind* able to define the humanistic content of a project, since the hierarchical architecture condenses the belief of a large number of users together with the careful consideration of many experts, in combination of all types of data and information, allowing the project HG to be derived.

Optimal Controller Design for Finite Word Length Implementation Using Genetic Learning Algorithm

Wen-Shyong Yu

Department of Electrical Engineering, Tatung Institute of Technology,
40 Chung-Shan North Rd. 3rd. Sec., Taipei, Taiwan 10451, Republic of China

ABSTRACT

In this paper, a linear quadratic gaussian (LQG) controller with genetic learning algorithm (GLA) is proposed to tackle the numerical errors due to the conversions of the A/D and D/A converters in a digital computer. This scheme can be directly used for the design of the ideal LQG and also is optimal in the presence of the numerical errors due to the finite word-length. By converting the stochastic problem to a deterministic game theoretic one, we find the estimation states using GLA and controller can minimize a suitable performance measure. The GLA, via reproduction, crossover, and mutation procedures, is used to tackle the signals from ADC to reduce the numerical errors and to obtain their optimal values.

INTRODUCTION

In digital control systems, it is necessary to quantize or round off the analog signals coming from the physical systems into the nearest digital level by A/D process and to the most likely analog signals by D/A process. Regardless of whether integer or fractional coding is utilized, the computation word provides a resolution of one part in 2^n where n is the number of bits, and the higher number of bits, n , the quantization error is more effectively decreased. Therefore, the numerical errors of the computation can drastically affect the performance of the control system. There exist realizations of a given controller transfer function yielding large effects from computational errors. Therefore, it is important to have a systematic way of reducing these effects [1-4,7]. But it was not until the last decade, that the problem was solved and not much concerned about by researches. In this paper, a linear quadratic Gaussian (LQG) controller with genetic learning algorithm (GLA) is proposed to tackle the numerical errors due to the conversions of the A/D and D/A converters in a digital computer. This control scheme can be directly used for the design of the ideal LQG and also is optimal in the presence of numerical errors due to the finite word length. The scheme is to convert the stochastic problem to a deterministic game theoretic one. We find the estimation states using GLA and design the controller so as to minimize a suitable performance measure. A necessary condition is developed for the existence of a saddle point solution.

ROUND-OFF ERROR AND LQG CONTROLLER DESIGN PROBLEM

Let the plant to be controlled be described as the time-invariant discrete-time single-input single-output (SISO) mathematical model:

$$\begin{cases} x_p(k+1) = A_p x_p(k) + \sum_{i=1}^q E_i x_p(k) \beta_i(k) + B_p(k) u(k) + H w_p(k) \\ y(k) = C_p^T x_p(k) + \sum_{i=1}^{\ell} F_i^T x_p(k) \psi_i(k) + G^T v_p(k) \end{cases} \quad k \in [0, N] \quad 1.$$

where $x_p(k) \in R^n$ is the state vector, $u(k)$ and $y(k)$ are the control vector and output, respectively, $\beta_i(k)$, $i = 1, 2, \dots, q$, $\psi_i(k)$, $i = 1, 2, \dots, \ell$, are uncorrelated white sequences, uncorrelated of each other, and $w_p(k)$, $v_p(k)$ are white sequences uncorrelated of each other with components uncorrelated

of $\beta_i(k)$ and $\psi_i(k)$; $A_p, E_i, i = 1, 2, \dots, q$, are $n \times n$ matrices, C_p and $F_i(k), i = 1, 2, \dots, \ell$, $n \times 1$ matrices, and H, G scalars. Let $E[\beta_i(j)\beta_i(k)] = \theta_i \delta_{jk}$, $E[\beta_i(j)\beta_m(k)] = 0$, $E[w_p(j)w_p(k)] = W\delta_{jk}$, $E[v_p(j)v_p(k)] = V\delta_{jk}$, $E[\psi_i(j)\psi_i(k)] = \phi_i \delta_{jk}$, and $E[\psi_i(j)\psi_m(k)] = 0$, for some $W > 0, V > 0$, and the white sequences have zero mean. We assume that A_p, B_p, C_p , and all the covariance matrices are known in advance.

Consider the following unbiased estimator and the LQG controller for (1):

$$\hat{x}_p(k+1) = A_p \hat{x}_p(k) + B_p u(k) + L(y(k) - C_p^T \hat{x}_p(k)) \quad 2.$$

$$u = -M \hat{x}_p(k) \quad 3.$$

where $\hat{x}_p(k)$ is the estimate of $x_p(k)$.

Define the estimation error by $\tilde{x}_p(k) = x_p(k) - \hat{x}_p(k)$ and the quantization errors by

$$Q[u(k)] = u(k) + e_u(k), \quad Q[\hat{x}_p(k)] = \hat{x}_p(k) + e_x(k) \quad 4.$$

Then we have

$$Q[\tilde{x}_p(k)] = Q[\hat{x}_p(k)] - x_p(k) = \hat{x}_p(k) - x_p(k) + e_x(k) = \tilde{x}_p(k) + e_x(k) \quad 5.$$

We shall ignore coefficient errors in (A_p, B_p, C_p, M) , since it is evident that controller structures that are good with respect to state quantization tend to also be good with respect to coefficient quantization. We then obtain the following closed loop system including finite word length effects:

$$\begin{cases} x_p(k+1) = A_p x_p(k) + \sum_{i=1}^q E_i x_p(k) \beta_i(k) + B_p u(k) + B_p e_u(k) + H w_p(k) \\ \tilde{x}(k+1) = (A_p - LC_p^T) \tilde{x}_p(k) + (A_p - LC_p^T) e_x(k) + H w_p(k) - LG v_p(k) \\ \quad + \sum_{i=1}^q E_i x_p(k) \beta_i(k) - L \sum_{i=1}^{\ell} F_i^T x_p(k) \psi_i(k) \\ u(k) = -M \hat{x}_p(k) - M e_u(k) \end{cases} \quad 6.$$

Define

$$x(k) = \begin{bmatrix} x_p(k) \\ \tilde{x}(k) \end{bmatrix}, A = \begin{bmatrix} A_p - B_p M & -B_p M \\ 0 & A_p - LC_p^T \end{bmatrix}, \bar{E}_i = \begin{bmatrix} E_i & 0 \\ E_i & 0 \end{bmatrix}, F_j = \begin{bmatrix} 0 & 0 \\ -LF_j^T & 0 \end{bmatrix},$$

$$\bar{H} = \begin{bmatrix} H & 0 \\ H & -LG \end{bmatrix}, w(k) = \begin{bmatrix} w_p(k) \\ v_p(k) \end{bmatrix}, e(k) = \begin{bmatrix} e_x(k) \\ e_u(k) \end{bmatrix}, A_e = \begin{bmatrix} 0 & -B_p M \\ A_p - LC_p^T & 0 \end{bmatrix}.$$

Then we have from Equation 6.

$$x(k+1) = Ax(k) + \sum_{i=1}^q \bar{E}_i x_p(k) \beta_i(k) + \bar{H} w(k) + A_e e(k) + \sum_{i=1}^{\ell} \bar{F}_i x_p(k) \psi_i(k) \quad 7.$$

The purpose of this paper is to find L and M such that the following cost function

$$J = E \left\{ \tilde{x}_p^T(N) Q_1 \tilde{x}_p(N) + x_p^T(N) Q_2 x_p(N) + \sum_{k=0}^{N-1} [x_p^T(k) Q_3 x_p(k) + u(k) Q_4 u(k)] \right\} + \sum_{i=1}^m \rho_i k_i^{-1} \quad 8.$$

is minimized where Q_1 , Q_2 , Q_3 , and Q_4 are suitable symmetric, positive definite matrices, ρ_i is the weighting factor, and $k_i = \frac{1}{12} 2^{-2n_i}$ where n_i is the word length used to store state variable $\hat{x}(k)$ in the digital computer.

Let the second moment of $x(k)$ be defined by

$$P(k) = E[x(k)x^T(k)] = \begin{bmatrix} p_{11}(k) & p_{12}(k) \\ p_{21}(k) & p_{22}(k) \end{bmatrix} \quad 9.$$

From Equation 7., we obtain

$$P(k+1) = AP(k)A^T + \sum_{i=1}^q \bar{E}_i P(k) \bar{E}_i^T \theta_i + \bar{H} \Lambda_1 \bar{H}^T + A_e \Lambda_2 A_e^T + \sum_{i=1}^{\ell} \bar{F}_i P(k) \bar{F}_i^T \phi_i \quad 10.$$

where $\Lambda_1 = \text{diag}\{W, V\}$ and $\Lambda_2 = \text{diag}\{K, R\}$ where $K_{ij} = k_i \delta_{ij}$, and $R = \frac{1}{12} 2^{-2\alpha} I$ is the

covariance of $e_u(k)$ for which α is the fractional part of the word length of the D/A converter. Let the Hamiltonian be defined by

$$H = \text{tr} \left[\left(AP(k)A^T + \sum_{i=1}^q \bar{E}_i P(k) \bar{E}_i^T \theta_i + \bar{H} \Lambda_1 \bar{H}^T + A_e \Lambda_2 A_e^T + \sum_{i=1}^{\ell} \bar{F}_i P(k) \bar{F}_i^T \phi_i \right) F(k-1) \right] \quad 11.$$

where $F(k-1) = \begin{bmatrix} f_{11}(k-1) & f_{12}(k-1) \\ f_{21}(k-1) & f_{22}(k-1) \end{bmatrix}$ is the costate matrix. Defining $F(k) = \partial H / \partial P(k)$, and

taking $\partial H / \partial M = \partial H / \partial L = 0$, we obtain

$$F(k) = A^T F(k-1)A + \sum_{i=1}^q \bar{E}_i F(k-1) \bar{E}_i^T \theta_i + \bar{H} \Lambda_1 \bar{H}^T + A_e \Lambda_2 A_e^T + \sum_{i=1}^{\ell} \bar{F}_i F(k-1) \bar{F}_i^T \phi_i \quad 12.$$

$$L = A_p p_{22} C_p^T \left(C_p p_{22} C_p^T + GVG^T + \sum_{i=1}^{\ell} \bar{F}_i p_{11} \bar{F}_i^T \phi_i \right)^{-1} \quad 13.$$

$$M = (Q_4 + B_p^T f_{11} B_p)^{-1} B_p^T f_{11} A_p \quad 14.$$

where $f_{12} = 0$ and $p_{12} = p_{22}$.

GENETIC LEARNING ALGORITHM

Given the I/O data set $X = \{x_1, x_2, \dots, x_n\}$, we can formulate a constrained optimization problem to let the signals from A/D converter approximately equal the true ones. The proposed GLA scheme explores a population of signals in parallel. Each signal in the population is coded as a string, and a collection of strings forms a generation. Except the input signal, the others are selected nearly from the input. These signals constitute the set X and are used to generate a more accurate signal by GLA in a digital computer. Given the n input/output data set and an initial population $G(0)$, the GLA generates a new generation $G(k+1)$ based on the previous generation $G(k)$ as follows:

- Step 1: $k = 0$;
 Step 2: Generate an initial population $G(0)$;
 Step 3: Evaluate $G(k)$;
 Step 4: If some termination conditions are met, go to Step 9;
 Step 5: Generate new generation $G(k+1)$ from $G(k)$;
 Step 6: Evaluate $G(k+1)$;
 Step 7: $k = k+1$;
 Step 8: Return to Step 3;
 Step 9: Stop.

The evaluation function that we use is given by

$$J_{GLA} = \frac{1}{1 + e^2} \quad 15.$$

where e^2 is the square error of the difference between output and reference input. In most cases we base the initialization procedure of the GLA on the initial selection of the mathematical model of the plant. The termination condition for the GLA occurs when the maximum generation number is reached.

Theorem 1: If $(L^*, M^*, \theta_i^*, V^*, W^*, \phi^*, \Lambda_1^*, \Lambda_2^*)$ is a saddle point (i.e. L^*, M^* are optimal solution and $\theta_i^*, V^*, W^*, \phi_i^*, \Lambda_1^*, \Lambda_2^*$ are all convex and compact set) for (10) and (12) to minimize (8), then

$$tr \left[\sum_{i=1}^{\ell} \bar{F}_i P \bar{F}_i^T \phi_i^* F(k-1) \right] \geq tr \left[\sum_{i=1}^{\ell} \bar{F}_i P \bar{F}_i^T \phi_i F(k-1) \right] \quad 16.$$

$$tr \left[\sum_{i=1}^{\ell} \bar{E}_i P \bar{E}_i^T \theta_i^* F(k-1) \right] \geq tr \left[\sum_{i=1}^{\ell} \bar{E}_i P \bar{E}_i^T \theta_i F(k-1) \right] \quad 17.$$

$$tr [H \Lambda_1^* H^T (f_{11}(k-1) + f_{22}(k-1))] \geq tr [H \Lambda_1 H^T (f_{11}(k-1) + f_{22}(k-1))] \quad 18.$$

$$tr [A_e \Lambda_2^* A_e^T (f_{11}(k-1) + f_{22}(k-1))] \geq tr [A_e \Lambda_2 A_e^T (f_{11}(k-1) + f_{22}(k-1))] \quad 19.$$

$$tr [L^* G V^* G^T L^* f_{22}(k-1)] \geq tr [L^* G V G^T L^* f_{22}(k-1)] \quad 20.$$

▽▽▽

Theorem 2: Suppose that a solution of $f_{11} > 0$, $f_{22} > 0$, $p_{11} > 0$, $p_{22} > 0$ exists for equations (10) and (12) for the point $(L, M, \theta_i, V, W, \phi, \Lambda_1, \Lambda_2)$, then the necessary conditions (15) -(19) of Theorem 1 are valid and the system is stabilized in the mean for the point (L, M) given by (13) and (14). Regardless of optimality, if a point $(L, M, \theta_i, V, W, \phi, \Lambda_1, \Lambda_2)$ guarantees such a solution, then the corresponding (L, M) stabilizes the system.

▽▽▽

EXAMPLE

Let

$$A = \begin{bmatrix} -7.020 \times 10^{-3} & 6.339 \times 10^{-1} & 5.180 \times 10^{-3} & -5.557 \times 10^{-1} & -6.112 \times 10^{-2} & 0 \\ -1.654 \times 10^{-2} & -3.889 \times 10^{-1} & 1.006 & 5.910 \times 10^{-3} & -4.632 \times 10^{-2} & 0 \\ 6.100 \times 10^{-4} & -3.521 \times 10^{-1} & -4.738 \times 10^{-1} & 0 & 1.783 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & -20 & 20 \\ 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$$B = [0 \ 0 \ 0 \ 0 \ 0 \ 30]^T, \quad C = [0 \ 0 \ 1 \ 0 \ 0 \ 0]$$

We assume that the continuous-time model (1) is digitally controlled using a pulse-amplitude-modulation signal with constant sampling frequency of 55Hz. We also assume that measurement of $y(t)$ is first subjected to an ideal anti-aliasing filter of bandwidth $2t_s^{-1}$ prior to sampling so that the discrete measurement covariance Λ_1 is given by

$$\Lambda_1 = \frac{10^{-6}}{t_s} = 5 \times 10^{-5}$$

Finally, the results are shown in Table 1.

Table 1. The performance of the system ($J \times 10^5$)

Number of bits	Non-GA design method	GA design method
12	9.6544	9.0010
10	9.8852	9.1101
8	9.9941	9.2102
6	12.3211	9.4011

CONCLUSION

In this paper, we have proposed a linear quadratic gaussian (LQG) controller using genetic learning algorithm (GLA) to tackle the numerical errors due to the conversions of the A/D and D/A converters in a digital computer. This scheme can be directly used for the design of the ideal LQG and also is optimal in the presence of the numerical errors due to the finite wordlength. We find the estimation states using GLA and controller can minimize a suitable performance measure. Since there are multiplicative noise terms, the separation principle is unfortunately not valid in this case and a set of necessary but not sufficient conditions are then proposed to stabilize the system.

REFERENCES

1. K Liu, R.E. Skelton, and K. Grigoriadis, 1992, Optimal controllers for finite word length implementation, IEEE Trans. on Automatic Contr., VOL. 37, NO. 9, pp. 1294-1304.
2. K. Kadiman and D. Williamson, 1989. Optimal finite word length linear quadratic regulation. IEEE Trans. on Automatic Contr., 34(12), 1218-1228.
3. P. Moroney, A. Willsley, and P. Houpt, 1983. Round-off noise and scaling in the digital implementation of control compensators. IEEE Trans. on Acoust. Speech, Signal Processing, 31, 1464-1477.
4. D. Williamson, 1985. Finite word length design of digital Kalman filters for state estimation. IEEE Trans. on Automatic Contr., 30, 30-39.
5. A. Sripad and D. Snyder, 1977. A necessary and sufficient condition for quantization error to be uniform and white, IEEE Trans. on Acoust. Speech, Signal Processing, 5, 442-448.
6. H. Kwakernaak and R. Silvan, 1989. Linear Optimal Control Systems, New York: Wiley.
7. D.E. Goldberg, Genetic Algorithm in Search, Optimization, and Machine Learning, Addison-Wesley Publishing Company, Inc..

ADAPTIVE FUZZY CONTROLLER FOR NONLINEAR UNCERTAIN SYSTEMS

Chiang-Cheng Chiang and Chih-Chien Hu

*Department of Electrical Engineering, Tatung Institute of Technology, 40,
Chung-Shan North Road, Sec. 3, Taipei, Taiwan, Republic of China*

ABSTRACT

In this paper, an adaptive fuzzy controller for single-input single-output (SISO) nonlinear systems with higher-order and unmatched uncertainties is provided. The design is based on the principle of sliding mode control, input-output (i/o) linearizing controller, and the approximation capability of adaptive fuzzy systems. Finally, the suggested controller is applied to the controlling of the Robot with Flexible Joint. Simulation results indicate its efficiency.

INTRODUCTION

Recently, much attention has been paid to the problem of designing a robust controller for a dynamical nonlinear system containing uncertainties and disturbances. Although the technique of the feedback linearization can achieve exact linearization of the closed-loop system, it requires the accurate mathematical model of the system because the differentiation procedure cannot be performed. In practice, however, many single-input single-output nonlinear systems are described by $\dot{x} = F(x, u)$, $y = h(x)$, where $x \in R^n$ is the state vector, $u \in R$ and $y \in R$ are the input and output of the system, respectively, and F and h are unknown nonlinear function. Thus, one may think to replace the F and h by fuzzy logic systems, and to develop an adaptive law to adjust the parameters of the fuzzy logic systems such that they will approximate F and h . Indeed, this is the basic idea of Sastry and Isidori [2] where the F and h are approximated by series expansions of known nonlinear basis functions. Therefore, another approach that approximates the nonlinear function in the final equation of the differentiation procedure by fuzzy logic systems and develops an adaptive law to adjust the parameters of the fuzzy logic systems is proposed by Wang [4]. Thus, the systematic design of adaptive fuzzy controllers using input-output linearization concept with mismatched uncertainties is conducted, and two main extensions is used: (1) variable structure control [5] (2) Lyapunov-based control [6, 8], to increase the robustness and to improve the performance of the controlled nonlinear system.

I/O LINEARIZATION OF NONLINEAR SYSTEMS WITH HIGHER-ORDER AND UNMATCHED UNCERTAINTIES

Consider the single-input single-output nonlinear system with system uncertainties:

$$\dot{x} = f(x) + g(x)u + \Theta(x); \quad y = h(x) \quad (1)$$

where $x \in R^n$ is the state, $u \in R$, $y \in R$ is the system input and output, respectively; f , $g \in R^n$ are known function and smooth vector fields. $\Theta \in R^n$ is the system uncertainty. $h(x)$ is a sufficiently smooth output function. The desired output trajectory is y_d . Assuming that the state coordinate transformation $(z, \eta) = T(x)$ is performed to the system (1), we can obtain the following form:

$$\dot{z}_i = z_i + \Delta\phi_i(z, \eta), \quad i = 1, 2, \dots, r-1; \quad \dot{z}_r = v + \Delta\phi_r(z, \eta) \quad (2)$$

$$\dot{\eta} = q(z, \eta) + \Delta\Omega(z, \eta) \quad (3)$$

and

$$u = \frac{1}{a(z, \eta)} [-b(z, \eta) + v] \quad (4)$$

where $a(z, \eta) = L_g L_f^{r-1} h(x) \circ T^{-1}(x)$; $b(z, \eta) = L_f h(x) \circ T^{-1}(x)$

$\Delta\phi_i = L_{\Theta} L_f^{i-1} h(x) \circ T^{-1}(z, \eta)$, $i = 1, 2, \dots, r$;

$\Delta\Omega = [L_{\Theta}\eta_1 \quad L_{\Theta}\eta_2 \quad \dots \quad L_{\Theta}\eta_{n-r}] \circ T^{-1}(z, \eta)$

Assumption 1: The zero dynamics is exponentially stable in the domain of definition, the function $q(z, \eta)$ is Lipschitz in z , and uniformly in η .

In the following assumptions, we can design a state feedback input-output linearizing controller (4) by the auxiliary control input v that can guarantee either exponential, or uniformly ultimately bounded stability of the nonlinear system (1) in the presence of uncertainties $\Theta(x)$ via sliding mode strategy with simple adaptive laws.

Assumption 2: The desired output trajectory y_d and its first r derivatives are uniformly bounded, that is

$$\|(y_d, y_d^{(1)}, \dots, y_d^{(r)})\| < B_d \quad (5)$$

where B_d is a positive constant.

Define the trajectory error to be

$$e_i = z_i - y_d^{(i-1)}, \quad i = 1, 2, \dots, r \quad (6)$$

Then the system (1) can be expressed as

$$\begin{aligned} \dot{e}_1 &= e_2 + \Delta\phi_1(e, \eta) \\ \dot{e}_2 &= e_3 + \Delta\phi_2(e, \eta) \\ &\vdots \end{aligned} \quad (7)$$

$$\dot{e}_{r-1} = e_r + \Delta\phi_{r-1}(e, \eta)$$

$$\dot{e}_r = v - y_d^{(r)} + \Delta\phi_r(e, \eta)$$

where $e = [e_1 \quad e_2 \quad \dots \quad e_r]^T$ is an error signal vector.

Assumption 3: The uncertainties are bounded by the polynomials combined with both $\|e\|^j$ and $\|\eta\|^k$, $j = 0, 1, \dots, N$, $k = 1, 2, \dots, M$. That is

$$\|\Delta\phi_i\| \leq \sum_{j=0}^N c_{1j}^i \|e\|^j + \sum_{k=1}^M c_{2k}^i \|\eta\|^k, \quad i = 1, 2, \dots, r \quad (8)$$

where c_{1j}^i and c_{2k}^i are positive constants, N and M are positive integers.

Assumption 4: The norm of the uncertainty vector $\Delta\Omega$ satisfies the following condition:

$$\|\Delta\Omega\| \leq L \quad \text{for all } z, \eta \in B_{(z, \eta)}. \quad (9)$$

We select a surface as follow:

$$S = e_r + a_1 e_{r-1} + \dots + a_{r-1} e_1, \quad a_j > 0 \text{ for } 1 \leq j \leq r-1 \quad (10)$$

The sliding surface is defined as $S = 0$, where a_j , $j = 1, 2, \dots, r-1$ are chosen so that all roots of polynomial $P(s) = s^{(r-1)} + a_1 s^{(r-2)} + \dots + a_{r-1}$ are in the open left half-plane. Then,

$$\dot{S} = v - y_d^{(r)} + [a_1 e_r + \dots + a_{r-1} e_2] + [\Delta\phi_r + a_1 \Delta\phi_{r-1} + \dots + a_{r-1} \Delta\phi_1] \quad (11)$$

and set $\Delta\Phi = [\Delta\phi_r + a_1 \Delta\phi_{r-1} + \dots + a_{r-1} \Delta\phi_1]$.

Consider the simple adaptive laws:

$$\dot{\hat{\Psi}}_{1j}(e) = q_{1j} |S| \|e\|^j, \quad j = 0, 1, 2, \dots, N \quad (12.a)$$

$$\dot{\hat{\Psi}}_{2k}(\eta) = q_{2k} |S| \|\eta\|^k, \quad k = 1, 2, \dots, M \quad (12.b)$$

where

$$\tilde{\Psi}_{1j}(e) = \bar{\Psi}_{1j}(e) - \sigma_{1j}, \quad \sigma_{1j} = \sum_{i=1}^r a_{r-i} c_{1j}^i, \quad j = 0, 1, 2, \dots, N \quad (12.c)$$

$$\tilde{\Psi}_{2k}(\eta) = \bar{\Psi}_{2k}(\eta) - \sigma_{2k}, \quad \sigma_{2k} = \sum_{i=1}^r a_{r-i} c_{2k}^i, \quad a_0 = 1, \quad k = 1, 2, \dots, M \quad (12.d)$$

are parameter adaptation errors and $q_{1j}, q_{2k} \in R$ are adaptation gains with positive values.

Then the state feedback control law:

$$u = \frac{1}{a(z, \eta)} [-b(z, \eta) + v]$$

where

$$v = y_d^{(r)} - \sum_{i=1}^{r-1} a_i e_{r-i+1} - \rho(e, \eta) \operatorname{sgn}(S) - pS, \quad p > 0 \quad (13)$$

$$\rho(e, \eta) = \sum_{j=0}^N \bar{\Psi}_{1j}(e) \|e\|^j + \sum_{k=1}^M \bar{\Psi}_{2k}(\eta) \|\eta\|^k \quad (14)$$

ADAPTIVE FUZZY CONTROLLER DESIGN

Assumption 5: There are the following linguistic descriptions about the unknown functions $b(z, \eta)$ and $a(z, \eta)$ (from human experts):

$$R_b^{(r)} : \text{IF } z_1 \text{ is } A_1^r \text{ and } \dots \text{ and } z_n \text{ is } A_n^r, \text{ THEN } b(z, \eta) \text{ is } C^r \quad (15)$$

$$R_a^{(s)} : \text{IF } z_1 \text{ is } B_1^s \text{ and } \dots \text{ and } z_n \text{ is } B_n^s, \text{ THEN } a(z, \eta) \text{ is } D^s \quad (16)$$

respectively, where A_i^r, B_i^s, C^r , and D^s are fuzzy sets in R , $r = 1, 2, \dots, L_b$ and $s = 1, 2, \dots, L_a$. Then we replace $b(z, \eta)$ and $a(z, \eta)$ in (4) by the fuzzy systems $\hat{b}(z|\theta_b)$ and $\hat{a}(z|\theta_a)$, respectively. The resulting control law can be rewritten as follows:

$$u_c = \frac{1}{\hat{a}(z|\theta_a)} [-\hat{b}(z|\theta_b) + v] \quad (17)$$

Then we consider a Lyapunov function:

$$V = \frac{1}{2} S^2 + \frac{1}{2} \sum_{j=0}^N q_{1j}^{-1} \tilde{\Psi}_{1j}^2(e) + \frac{1}{2} \sum_{k=1}^M q_{2k}^{-1} \tilde{\Psi}_{2k}^2(\eta) \quad (18)$$

According to (8) and (14)

$$\dot{V} \leq -pS^2 + [b(z, \eta) - \hat{b}(z|\theta_b) + a(z, \eta)u_c - \hat{a}(z|\theta_a)u_c]S \quad (19)$$

From (19), we see that it is very difficult to design the u_c such that the last term of (19) is less than zero. Thus, to solve this problem we append another control term u_s that is called the supervisory control [4]. Then, the resulting adaptive fuzzy control will be chosen as

$$u = u_c + u_s. \quad (20)$$

Assumption 6: We can determine functions $b^U(z, \eta)$, $a^U(z, \eta)$, and $a_L(z, \eta)$ such that $|b(z, \eta)| \leq b^U(z, \eta)$ and $a_L(z, \eta) \leq a(z, \eta) \leq a^U(z, \eta)$ for $z \in U_c$, where $b^U(z, \eta) < \infty$, $a^U(z, \eta) < \infty$, and $a_L(z, \eta) > 0$ for $z \in U_c$.

We choose u_s as

$$u_s = -I^* \operatorname{sgn}(S) \frac{1}{a_L(z, \eta)} [|\hat{b}(z|\theta_b)| + |b^U(z, \eta)| + |\hat{a}(z|\theta_a)u_c| + |a^U(z, \eta)u_c|] \quad (21)$$

In the following, we develop an adaptive law to adjust the parameters in the fuzzy systems for the purpose of forcing the tracking error to converge to zero.

First, define

$$\theta_b^* = \arg \min_{\theta_b \in \Omega_b} [\sup_{z \in U_c} |\hat{b}(z|\theta_b) - b(z, \eta)|] \quad (22)$$

$$\theta_a^* = \arg \min_{\theta_a \in \Omega_a} [\sup_{z \in U_c} |\hat{a}(z|\theta_a) - a(z, \eta)|] \quad (23)$$

where Ω_b and Ω_a are constraint sets for θ_b and θ_a , specified by the designer. For Ω_b , we require that θ_b is bounded; that is,

$$\Omega_b = \{\theta_b : |\theta_b| \leq M_b\} \quad (24)$$

M_b is positive constant specified by the designer. For Ω_a , in addition to the constraints similar to (24), we also require that $\hat{a}(z|\theta_a)$ must be positive. Then, we have

$$\Omega_a = \{\theta_a : |\theta_a| \leq M_a, \bar{y}^l \geq \varepsilon\} \quad (25)$$

where M_a , ε are positive constants specified by the designer. Define the minimum approximation error

$$w = (\hat{b}(z|\theta_b) - b(z, \eta)) + (\hat{a}(z|\theta_a) - a(z, \eta))u_c \quad (26)$$

Then the error equation can be rewritten as

$$\begin{aligned} \dot{e}_r = & \phi_b^T \xi_b(z) + \phi_a^T \xi_a(z)u_c + w + \Delta\phi_r(z, \eta) \\ & - \sum_{i=1}^{r-1} a_i e_{r-i+1} - \rho(e, \eta) \operatorname{sgn}(S) - pS + a(z, \eta)u_s \end{aligned} \quad (27)$$

where $\phi_b = \theta_b^* - \theta_b$, $\phi_a = \theta_a^* - \theta_a$, and $\xi_b(z)$ and $\xi_a(z)$ are the fuzzy basis functions. Then we consider a Lyapunov function as:

$$V = \frac{1}{2}S^2 + \frac{1}{2r_1}\phi_b^T \dot{\phi}_b + \frac{1}{2r_2}\phi_a^T \dot{\phi}_a + \frac{1}{2}\sum_{j=0}^N q_{1j}^{-1} \tilde{\Psi}_{1j}^2(e) + \frac{1}{2}\sum_{k=1}^M q_{2k}^{-1} \tilde{\Psi}_{2k}^2(\eta)$$

Then we obtain

$$\begin{aligned} \dot{V} \leq & \phi_b^T \xi_b(z)S + \phi_a^T \xi_a(z)u_c S + wS - pS^2 + a(z, \eta)u_s S + \frac{1}{r_1}\phi_b^T \dot{\phi}_b + \frac{1}{r_2}\phi_a^T \dot{\phi}_a \\ = & \frac{1}{r_1}\phi_b^T [\dot{\phi}_b + r_1 \xi_b(z)S] + \frac{1}{r_2}\phi_a^T [\dot{\phi}_a + r_2 \xi_a(z)u_c S] - pS^2 + wS + a(z, \eta)u_s S \end{aligned} \quad (28)$$

where $\dot{\phi}_b = \dot{\theta}_b$ and $\dot{\phi}_a = \dot{\theta}_a$. From (21) and $a_L(z, \eta) > 0$, we have that $a(z, \eta)u_s S \leq 0$.

If we choose the adaptive law

$$\dot{\theta}_b = -r_1 \xi_b(z)S \quad (29)$$

$$\dot{\theta}_a = -r_2 \xi_a(z)u_s S \quad (30)$$

then from (28), we have

$$\dot{V} \leq -pS^2 + wS. \quad (31)$$

Since w is the minimum approximation error, we have $\dot{V} \leq 0$.

For alleviating the chattering phenomenon caused by the sign function, we can use the fuzzy controller output u_f to replace the sign function. The fuzzy control rules which determine u_f are defined as follow:

$$R_j : \text{IF } S \text{ is } \tilde{F}_j, \text{ THEN } u_f \text{ is } \tilde{G}_j, \quad j = -2, -1, 0, 1, 2 \quad (32)$$

where \tilde{F}_j and \tilde{G}_j are the membership functions of the fuzzy sets. Adopting the max-product compositional rule of inference and the method of the singleton fuzzification [3] and using the method of the center of gravity defuzzification. We can obtain the fuzzy controller output u_f to replace the sign function.

AN EXAMPLE AND SIMULATION RESULTS

The system can be modeled by the two equations [1, 7]:

$$I\ddot{q}_1 + mgl \sin q_1 + k(q_1 - q_2) = 0 \quad (33)$$

$$J\ddot{q}_2 - k(q_1 - q_2) = u$$

The system dynamics in a state-space representation with the state vector is chosen as

$x = [x_1 \ x_2 \ x_3 \ x_4]^T = [q_1 \ \dot{q}_1 \ q_2 \ \dot{q}_2]^T$ and $y = h(x) = x_1 = q_1$.

Consider the system including exogenous noise and internal uncertainties, that is, $\Theta(x) = [\mu_1 x_1 \ \mu_2 x_1 x_2 \ \mu_3 x_3 \ \mu_4 x_2^2]^T$. To obtain numerical results, we set $mgl = 10$,

$k = 100$, $I = 100$, $J = 10$. The desired output trajectory is given as $y_d = \sin \frac{2\pi}{5} t$ and

$e_i = z_i - y_d^{(i-1)}$, $i = 1, 2, 3, 4$. $b^U(z, \eta) = 17$, $a_L(z, \eta) = 1$, and $a^U(z, \eta) = 1$ is known from Assumption 6.

The sliding surface is defined as follows:

$$S = e_4 + 6e_3 + 11e_2 + 6e_1.$$

The bound polynomials of the uncertainties in Assumption 4 is evaluated as

$$\|\Delta\phi_i\| \leq \sum_{j=0}^2 \bar{c}_j^i \|e\|^j, \quad i = 1, 2, 3, 4$$

We use simple extensive adaptive laws to estimate the upper bound of each order as:

$$\bar{c}_j(e) = \bar{c}_{0j} + q_j \int_0^t |S| \|e\|^j dt, \quad j = 0, 1, 2$$

where \bar{c}_{0j} are initial values. By choosing appropriate $\{\bar{c}_{0j}\}$ and $\{q_j\}$, we can adjust the rate of parameter adaptation. Here we select $q_0 = q_1 = q_2 = 3$ and the sampling time $\Delta t = 0.002$ sec.

Then the switching feedback gain can be select as $\rho(e) = \sum_{j=0}^2 \bar{c}_j(e) \|e\|^j$ and $p = 3$. The simulation results are shown in Figs. 1-2.

CONCLUSIONS

The problem of adaptive fuzzy control for a class of nonlinear system with higher-order and unmatched uncertainties was studied in this paper. We have some attractive characteristics in the proposed method: (1) do not require an accurate mathematical model of the system when the system is under control, (2) use fuzzy IF-THEN rules describing the system directly into the controller, (3) use fuzzy control output to replace the sign function. The proposed adaptive fuzzy controller is illustrated by the Robot with Flexible Joint, and the simulation results show the efficiency of the proposed method.

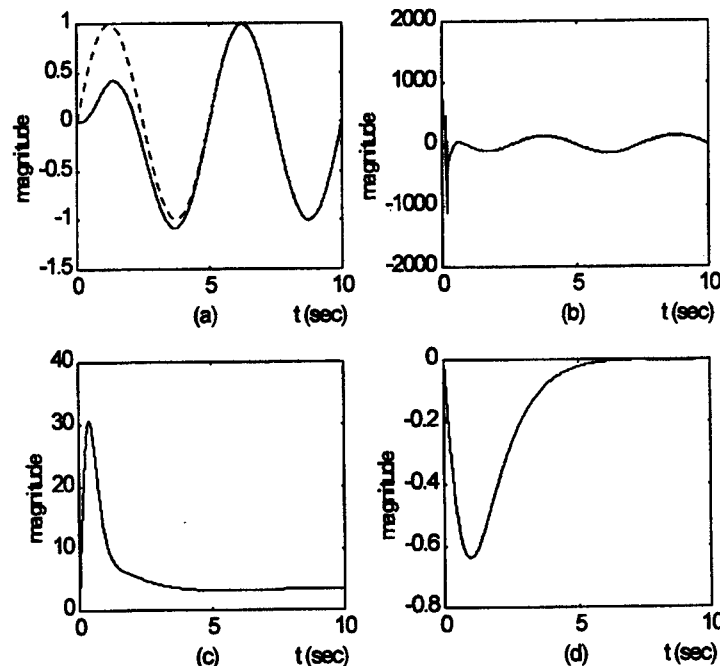


Fig. 1. System responses: (a) output $y(t)$; (b) input $u(t)$; (c) switching gain $\rho(t)$; (d) tracking error $e_1(t)$.

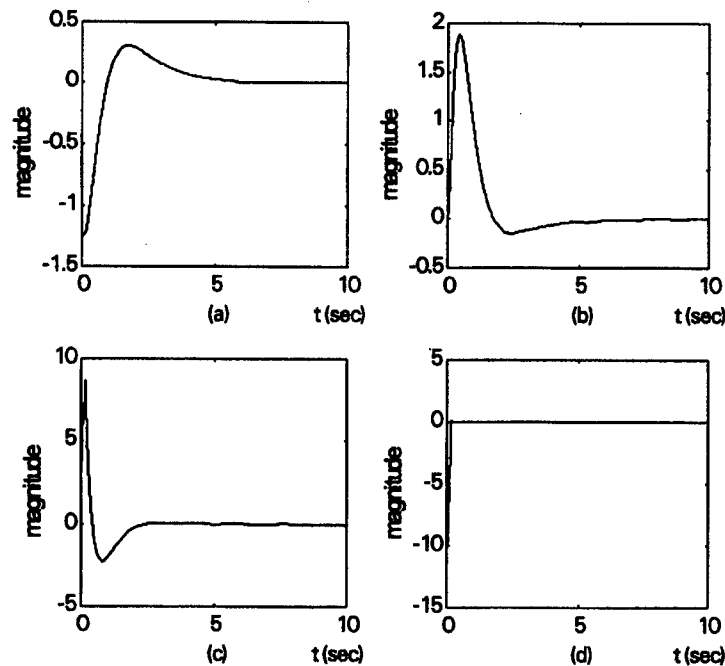


Fig. 2. System responses: (a) error $e_2(t)$; (b) error $e_3(t)$; (c) error $e_4(t)$; (d) sliding surface $S(t)$.

REFERENCES

- [1] A. Isidori, 1989. *Nonlinear control systems: an introduction*, 2nd ed.. New York: Springer-Verlag.
- [2] S. Sastry, and A. Isidori, 1989. Adaptive Control of linearizable systems. *IEEE Trans. On Automatic Control*, 34, 1123-1131.
- [3] L. X. Wang, 1993. Stable adaptive fuzzy control of nonlinear system. *IEEE Trans. Fuzzy Systems*, 1, 146-155.
- [4] L. X. Wang, 1994. *Adaptive Fuzzy Systems and Control Design and Stability Analysis*. Englewood Cliffs, NJ: Prentice-Hall, Inc..
- [5] T. L. Liao, L. C. Fu, and C. F. Hsu, 1992. Output tracking control of nonlinear systems with mismatched uncertainties. *System and Control Letters*, 18, 39-47.
- [6] K. Y. Lian, L. C. Fu, and T. L. Liao, 1993. Robust output tracking nonlinear systems with weakly non-minimum phase. *Int. J. Contr.*, 58, 301-316.
- [7] J. J. E. Slotine and W. Li, 1991. *Applied nonlinear control*. Englewood Cliffs, NJ: Prentice-Hall, Inc..
- [8] J. J. E. Slotine and J. K. Hedrick, 1993. Robust input-output feedback linearization. *Int. J. Contr.*, 57, 1133-1139.
- [9] S. Sastry, and M. Bodson, 1989. *Adaptive Control: Stability, Convergence, and Robustness*. Englewood Cliffs, NJ: Prentice-Hall, Inc..
- [10] T. P. Zhang and C. B. Feng, 1997. Decentralized adaptive fuzzy control for large-scale nonlinear systems. *Fuzzy Sets and Systems*, 92, 61-70.

Hybrid Modeling
(abstracts and viewgraphs)

Holistic Strategies for Designing Multistage Material Processes

W. Garth Frazier, E. Medina

Materials Process Design Branch, Materials Directorate
Air Force Research Laboratory
Wright-Patterson Air Force Base, Ohio
Email: fraziewg@ml.wpafb.af.mil

An algorithm is presented for using dynamic state-variable models of microstructure evolution formulated from first principles and empirical data in combination with fuzzy design rules and simplified analytical models of the effect that product shape has on achieving desirable final products.

Previously existing algorithms and software solved for the constrained, optimal conditions by combining the primal optimisation criterion and constraints into a single optimality criterion using a weighting technique. While often effective, this approach had the disadvantages of requiring the user to estimate relative weights and causing the algorithm to have slow rates of convergence. Also, prior software did not provide for any support of the effects of final product shape. Through application of further mathematical analysis, fuzzy design rules that account for shape effects and numerical techniques, the following new capabilities have been added.

- explicit handling of constraints
(no need for the user to choose weights)
- fuzzification of the influence of shape and the design
- design of initial process conditions
(temperature, grain size, etc.)
- design of processing time
(processes can be optimized for minimum processing time)
- multi-stage dynamics
(sequential thermo-mechanical processes can be designed for overall optimization, not just one stage at a time)

The presentation will include an explanation of model formulation, design strategies, numerical algorithms, and a example application and results.

References

- W.G. Frazier, et. al., "Application of Control Theory Principles to the Optimization of Grain Size During Hot Extrusion," *Materials Science and Technology*, 14, January, 1998, 25-31
- J.C. Malas, W.G. Frazier, S. Venugopal, E.A. Medina, S. Medeiros, R. Srinivasan, R.D. Irwin, W.M. Mullins, and A. Chaudhary, "Optimization of Microstructure Development During Hot Working Using Control Theory," *Metallurgical and Materials Transactions A*, 28A(9), September, 1997, 1921-1930
- S. Venugopal, E.A. Medina, J.C. Malas, S. Medeiros, W.G. Frazier, W.M. Mullins, and R. Srinivasan, "Optimization of Microstructure During Deformation Processing Using Control Theory Principles," *Scripta Materialia*, 36(3), 1997, 347-353
- E.A. Medina, S. Venugopal, W.G. Frazier, S. Medeiros, W.M. Mullins, A. Chaudhary, R.D. Irwin, R. Srinivasan, and J.C. Malas, "Optimization of Microstructure Development: Application to Hot Metal Extrusion," *Journal of Materials Engineering and Performance*, 5(6), 1996, 743-752



HOLISTIC STRATEGIES FOR DESIGNING MULTISTAGE MATERIALS PROCESSES

W. G. Frazier

Air Force Research Laboratory

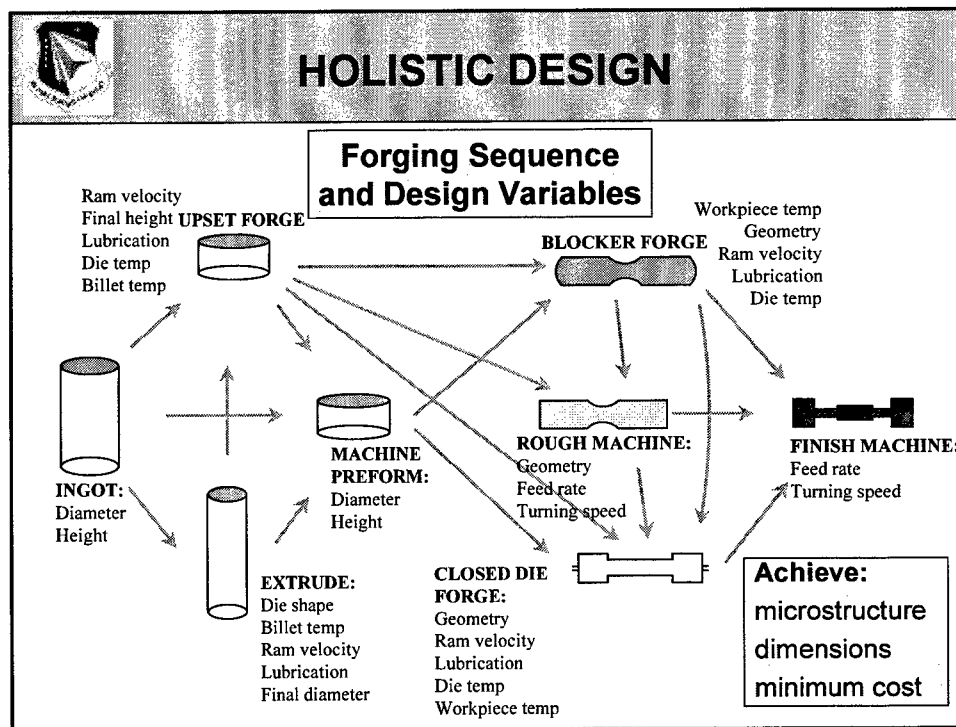
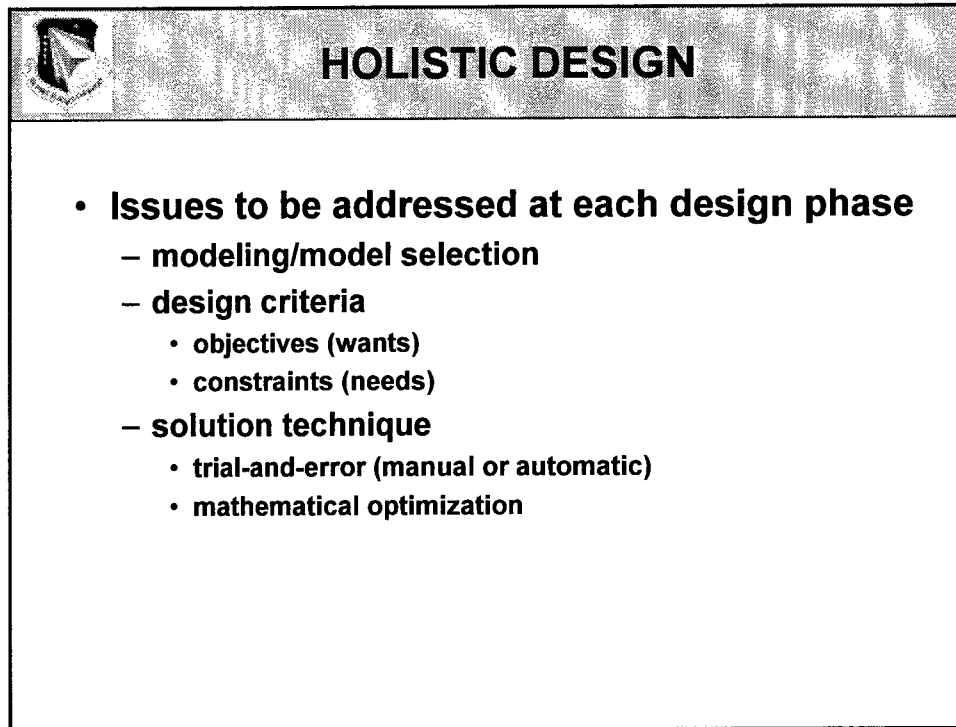
E. A. Medina

Austral Engineering and Software, Inc.



HOLISTIC DESIGN

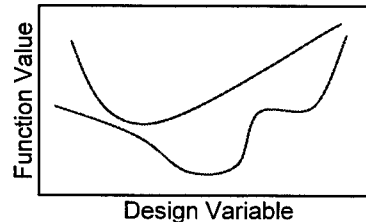
- **Why?**
 - Affordability goals require more than just “tuning” existing process sequences. New sequences and new processes must be competed.
- **Aspects**
 - when in the design process and where in the modeling structure do we assume independence?
 - breadth (the whole sequence should be considered at once and include coupling at least in the early design phases)
 - depth (effects on different scales should be considered at least in the final design phases)





MODELING FOR HOLISTIC DESIGN

- **Different degrees of model accuracy**
 - Tradeoffs among understanding, ease-of-use, functional requirements, impact on design phase
- **Features**
 - discrete, continuous
 - atomistic, continuum
 - static, dynamic
 - lumped, distributed
- **Sources**
 - first principles, analytical & numerical (science)
 - empirical (system identification)
 - heuristic (experiential)



DIFFICULTIES USING HYBRID MODELS

- **Integration of models of the same phenomenon**
 - FEA-based model vs. ERA-based model of a flexible structure
- **Embedding of types within a model**
 - FEA model of plastic deformation and CA model of microstructure evolution
- **Inconsistencies of accuracy/representation from process model to process model**
 - feature-based model of deformation and point-wise definition of thermal effects
 - feature-based model of heat transfer and point-wise definition of degree of macrosegregation

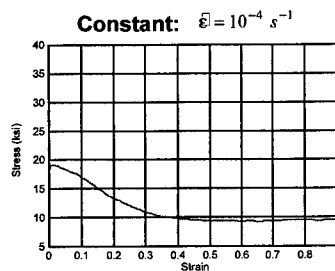


HYBRID MODELING AND HOLISTIC DESIGN

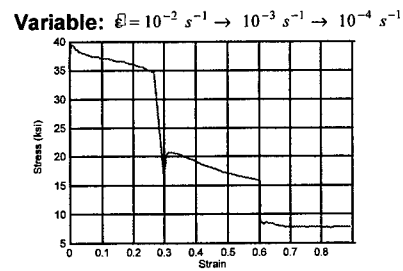
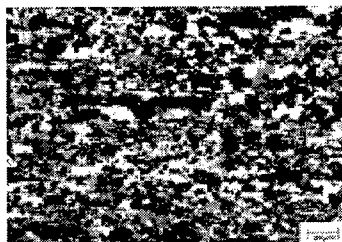
- Many models have limited domains of validity, thereby restricting design alternatives
 - has also lead automatic process control design to focus on regulation of process inputs, and not dynamical tracking and optimal trajectory correction of process “path”
- Need to open up degrees of freedom in the design space for alternative processing paths through dynamical model development
 - provides a means to optimal sequence selection and time-varying process input selection
- Material *models* tend to be process dependent and process objective driven



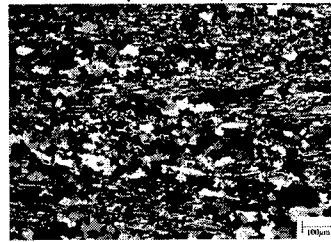
DYNAMICAL PROCESSING



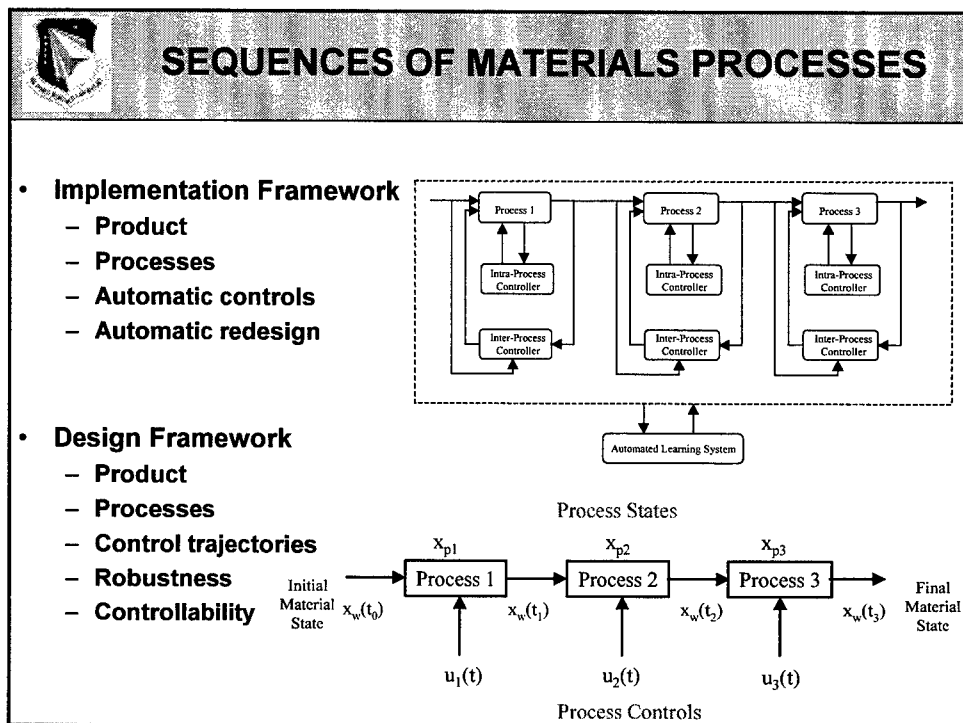
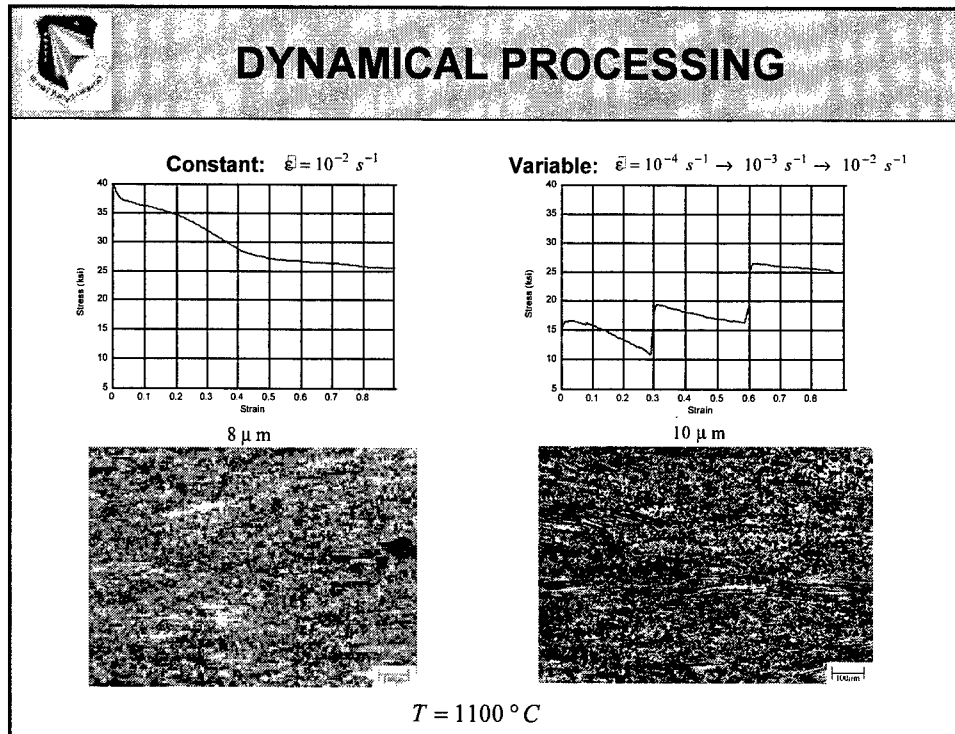
18 μm

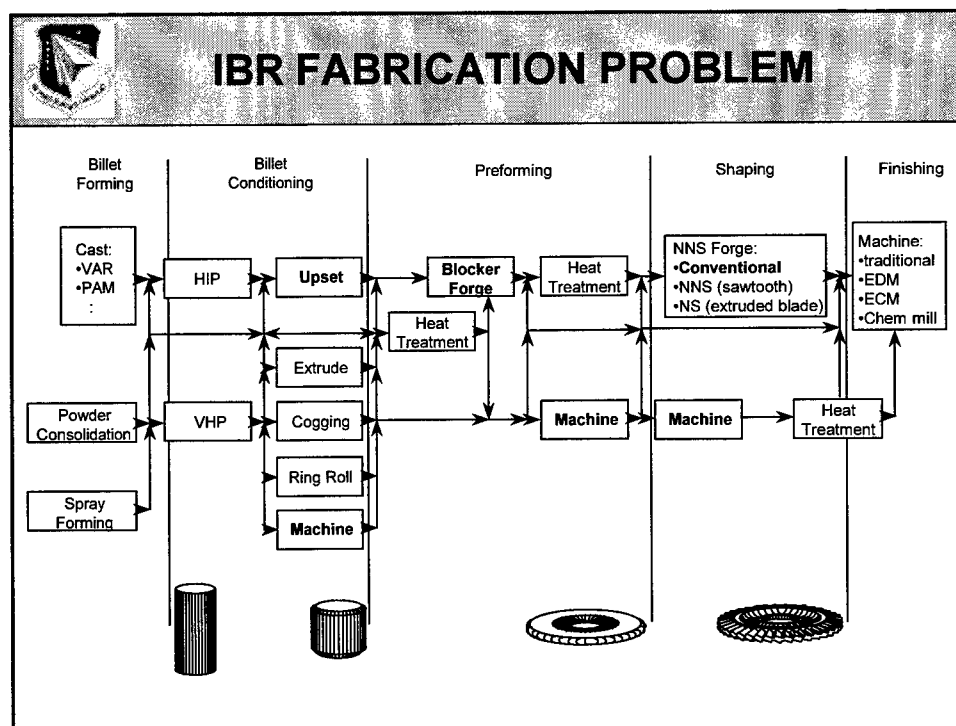
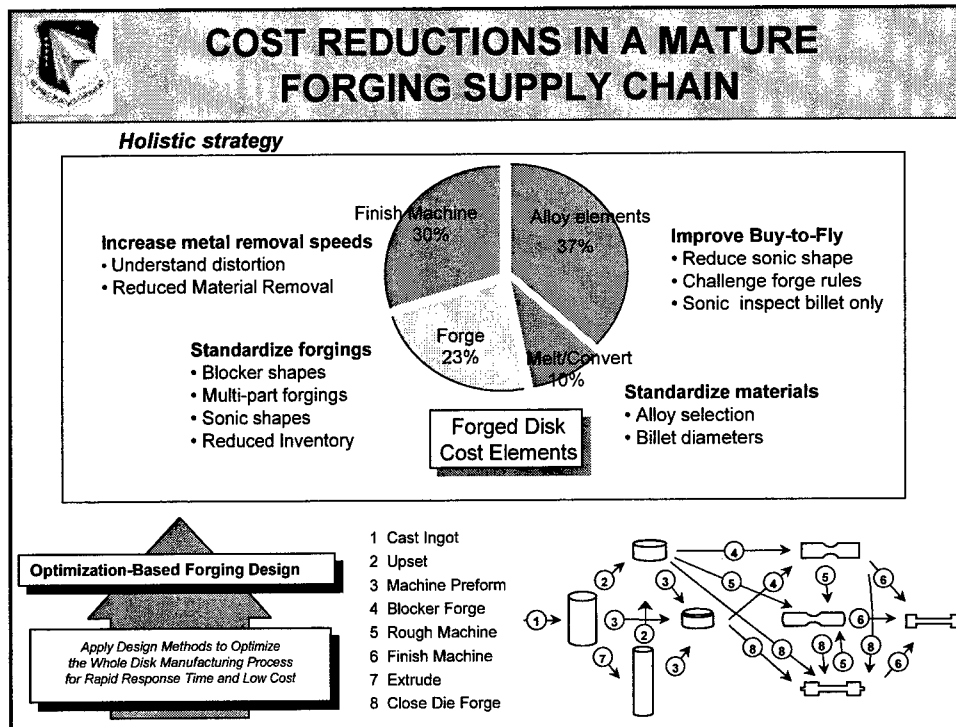


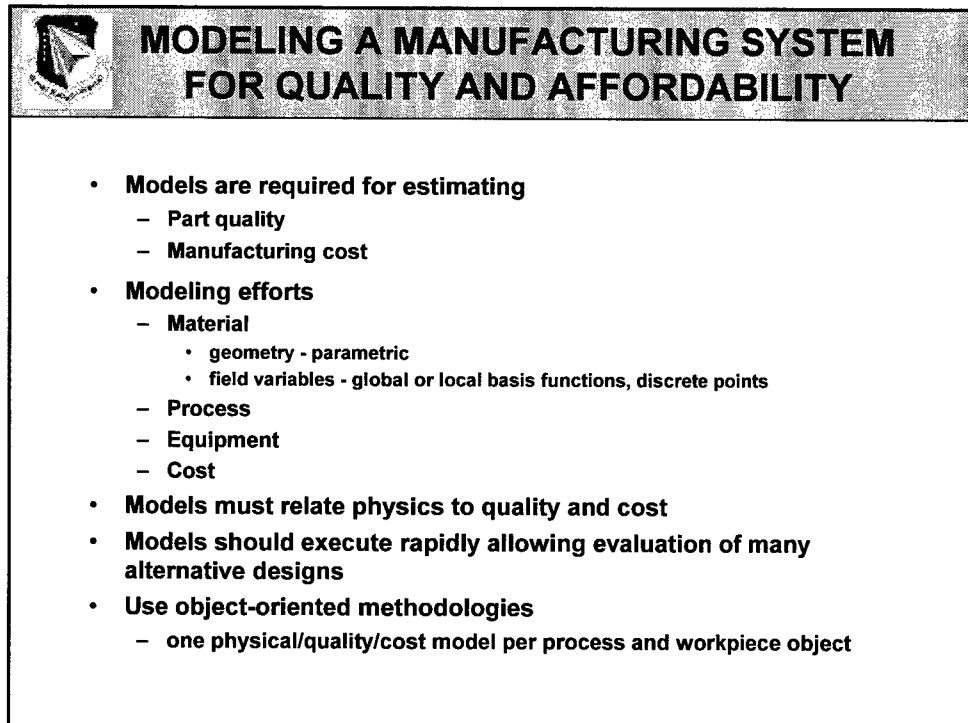
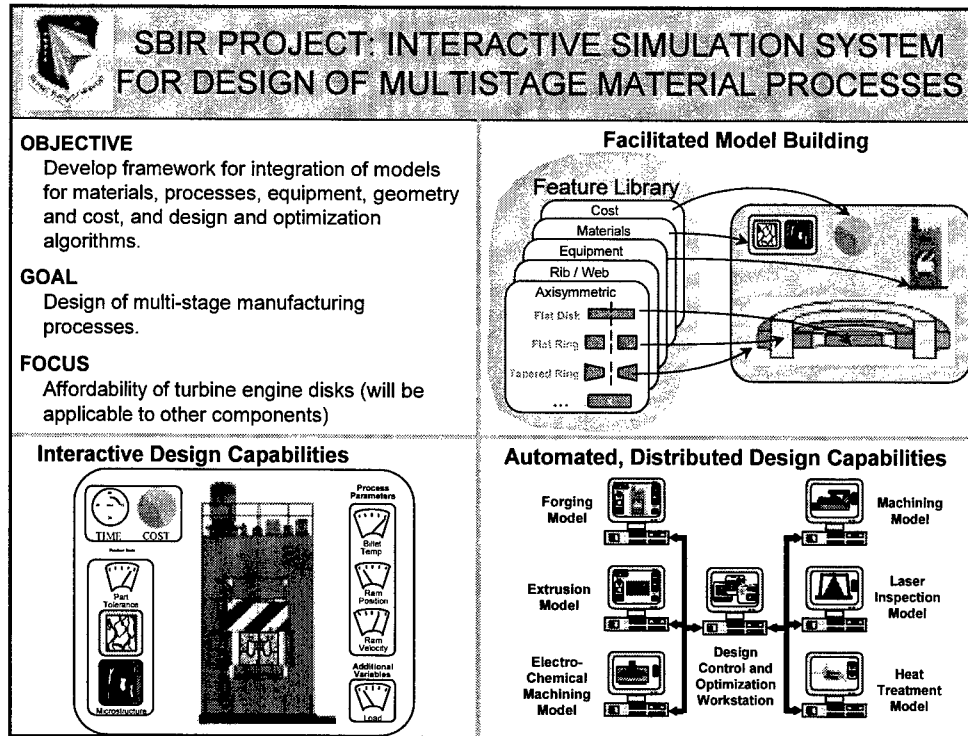
30 μm and 14 μm

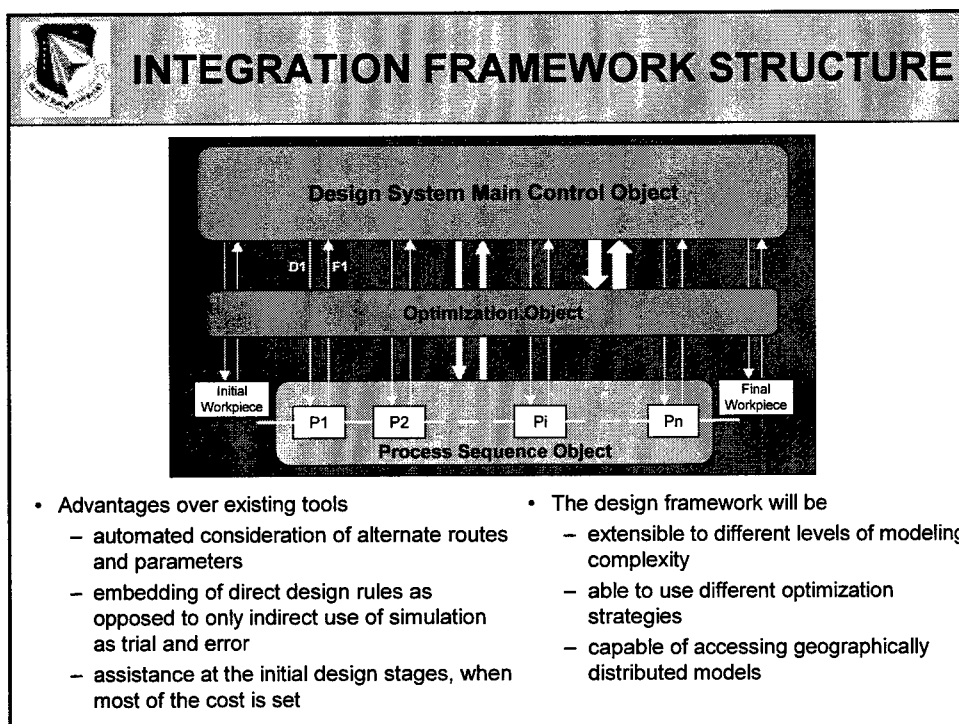
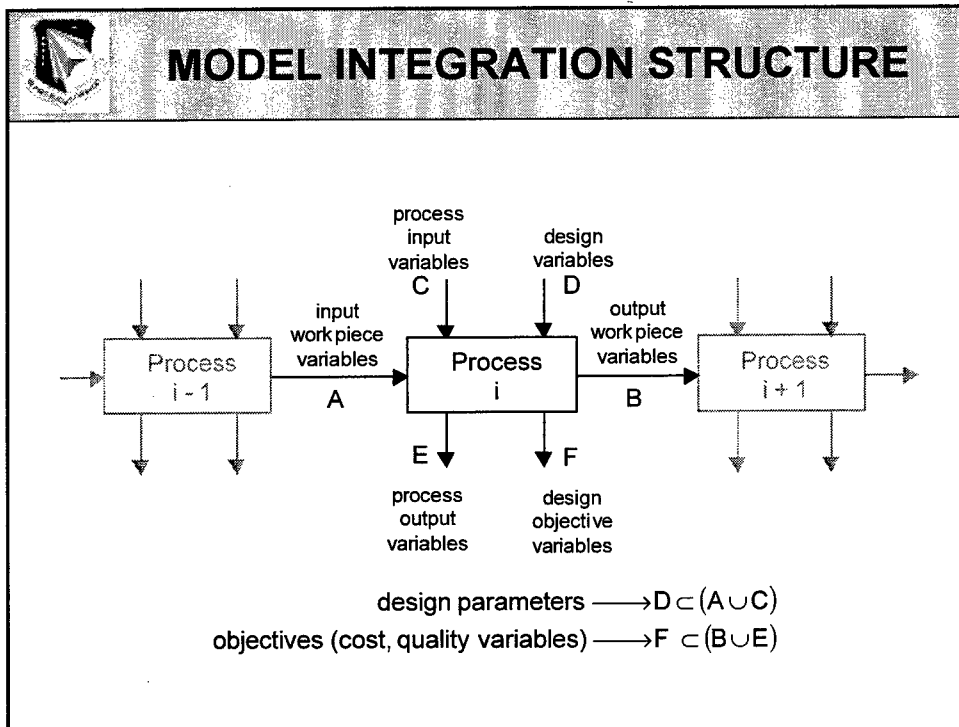


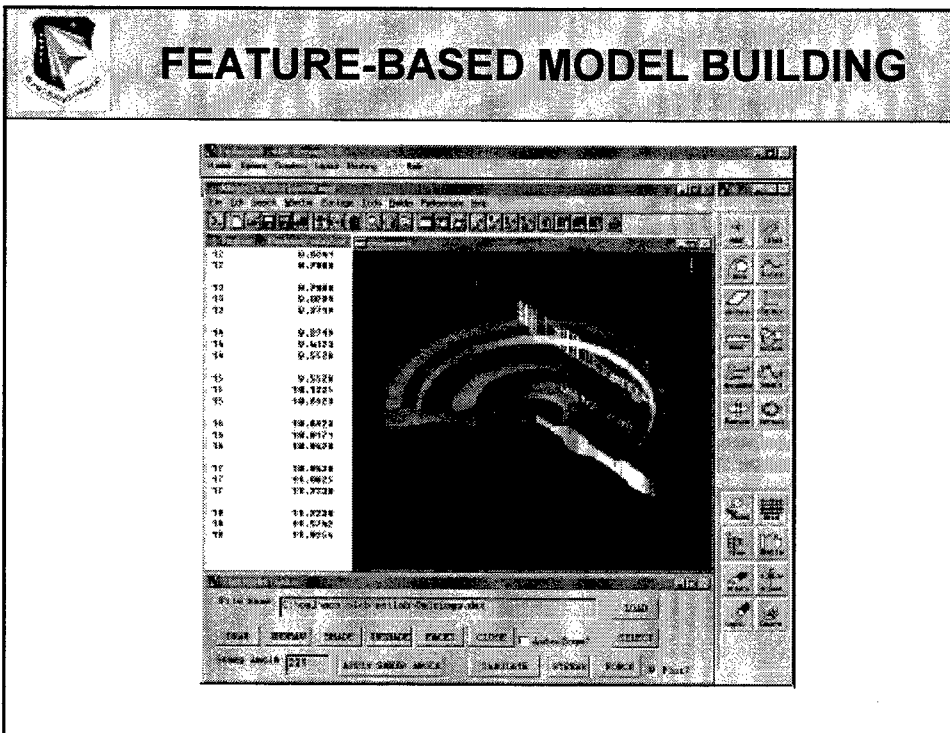
$T = 1100^{\circ}C$











DESIGN STRATEGIES

- **Gradient-based optimization**
 - SQP & related techniques
 - best for local optimization
- **Non-gradient-based optimization**
 - global optimization
- **Optimization using surrogates**
- **Ideal design**



CONTINUING ISSUES

- **Finding the correct model accuracy**
- **Integration framework needs to support all levels of model complexity**
- **Definition of design objectives and constraints**
- **Identification of decoupling strategies**

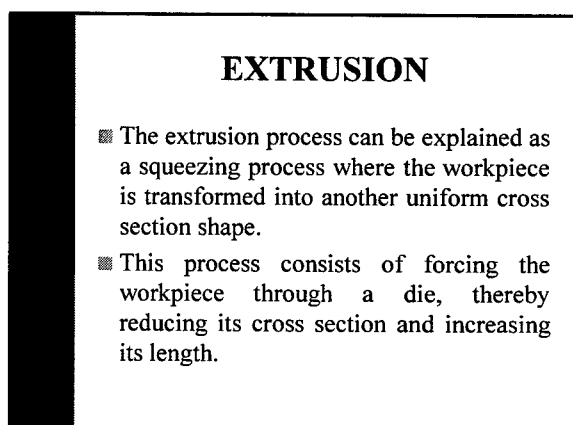
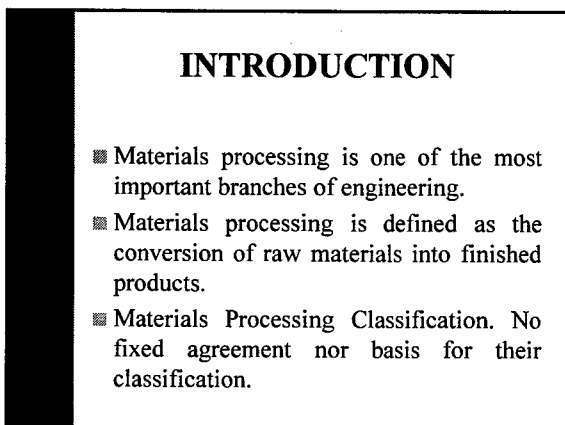
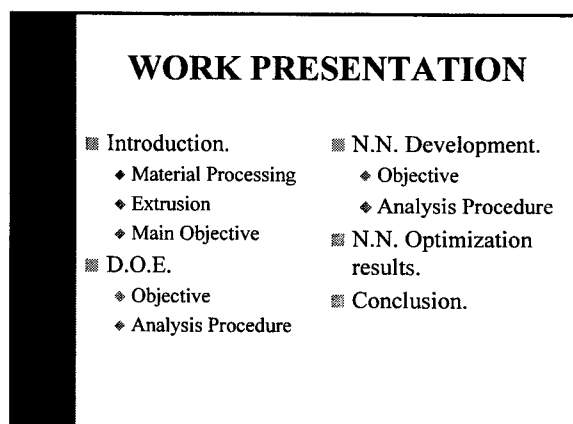
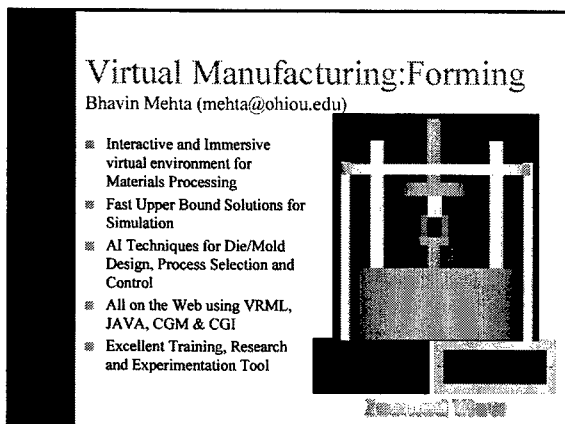
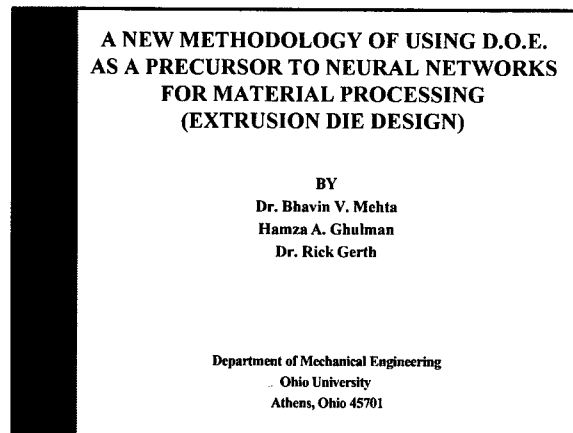
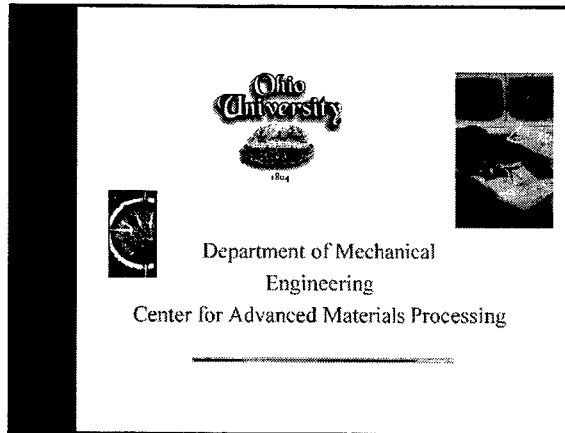
A New Methodology of Using Design of Experiments as a Precursor to Neural Networks for Material Processing: Extrusion Die Design

Bhavin V. Mehta*, Hamza Ghulman*, Rick Gerth**

* Department of Mechanical Engineering,
Ohio University, Athens, Ohio 45701, USA

** Department of Industrial and Manufacturing Engineering,
Ohio University, Athens, Ohio 45701, USA

Extrusion die design and making is an art and a science. In present day extrusions using composites, polymers, and other new alloys, the product geometries are extremely complicated. The flow analysis inside an extrusion die using Finite Element Analysis (FEA) is tedious and time consuming. To optimize the design of a die one needs to perform hundreds of runs, requiring several weeks or months of computer time. In the past researchers have used Neural Networks (NN) to optimize the design and predict flow patterns for newly designed dies of similar geometries. But, even for NN it has been proven that one needs a few thousand runs to train a network and accurately predict the flow. This paper shows a new methodology of using Design of Experiments (DOE) as a precursor to identify the importance of some variables and thus reduce the data set needed for training a NN. Based on the DOE results, a neural network training set is generated with more variations for the most significant inputs. A comparison of design using only NN versus using DOE and then NN is shown. The results indicate a significant reduction in the size of the training set, the time required for training and improvement in accuracy of the predicted results. To reduce the analysis time, a newly developed upper bound technique was used for generating the training set. The DOE model is extremely fast and can be used for real time (on-line) control of the process.



MAIN OBJECTIVE

- The objective of the project is to develop a new methodology that combines D.O.E. with N.N. to understand and analyze and predict the behavior of any process.
- Our sample case here is to predict the streamlined velocity fields of the extrusion process for the straight converging dies.

Design of Experiments "Objective"

- The main goal was to identify the significance of each parameter (variable) and its contribution to the model. Based on the results, one can decide on the training set for neural networks.

- The five main factors analyzed are :
Length of the die
Diameter ratio of the entry to exit
Initial velocity
n in the material properties equation
m for the friction coefficient
- For this kind of models, it is advisable to do the analysis using the Fractional Factorial (FF). The main advantage of this method is the reduction in the number of experiment.

- The main advantage of the FF is to determine the significant few from the insignificant.
- The analysis was performed for the model of 2^{5-1} with resolution V (five), which means that the model has 5 variables each one has two levels and one design generator I = ABCDE.
- The next analysis was carried out for the model of 2^{5-2} with resolution III and design generators I=ABD=ACE=BCDE.

**Data Distribution Sheet for Case 2^{5-1}
(E= ABCD)**

A	B	C	D	E
Length	Init. Vel.	Dia.Ratio	n (Mat.)	m (Fric.)
-1	-1	-1	-1	1
1	-1	-1	-1	-1
-1	1	-1	-1	-1
1	1	-1	-1	1
-1	-1	1	-1	-1
1	-1	1	-1	1
-1	1	1	-1	-1
1	1	1	-1	1
-1	-1	-1	1	-1
1	-1	-1	1	1
-1	1	-1	1	-1
1	1	-1	1	1
-1	-1	1	1	-1
1	-1	1	1	1
-1	1	1	1	-1
1	1	1	1	1

**Data Collection Sheet for Case 2^{5-1}
"variation due to n (material property)"**

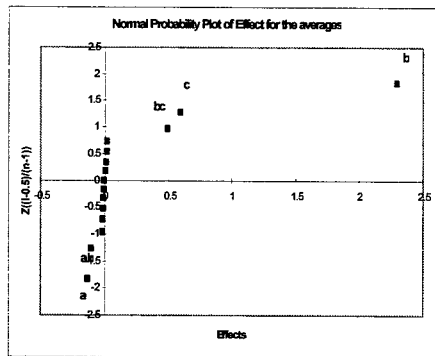
A	B	C	D	E	-0.01	+0.01	Average
Length	Init. Vel.	Dia.Ratio	n (Mat.)	m (Fric.)	"set I"	"set I"	
0.5	0.1	2	0.1	0.9	0.2173	0.2166	0.2169
8	0.1	2	0.1	0	0.1898	0.1896	0.1897
0.5	1	2	0.1	0	2.1115	2.1119	2.1117
8	1	2	0.1	0.9	1.9545	1.9522	1.9533
0.5	0.1	5	0.1	0	0.3192	0.3196	0.3194
8	0.1	5	0.1	0.9	0.2980	0.2976	0.2978
0.5	1	5	0.1	0.9	3.2251	3.2287	3.2269
8	1	5	0.1	0	2.9528	2.9515	2.9522
0.5	0.1	2	0.4	0	0.2118	0.2118	0.2118
8	0.1	2	0.4	0.9	0.1932	0.1930	0.1931
0.5	1	2	0.4	0.9	2.1169	2.1124	2.1146
8	1	2	0.4	0	1.8713	1.8696	1.8704
0.5	0.1	5	0.4	0.9	0.3264	0.3268	0.3266
8	0.1	5	0.4	0	0.2926	0.2926	0.2926
0.5	1	5	0.4	0	3.2494	3.2534	3.2514
8	1	5	0.4	0.9	2.9623	2.9614	2.9618

Data Analysis for Case 2⁵⁻¹

	Average	1	2	3	4	Effect	SS
(I)	0.2169	0.4066	4.4716	11.2679	22.4902	2.8113	***
A	0.1897	0.0650	6.7962	11.2223	-1.0683	-0.1335	0.0713
B	2.1117	0.6172	4.3899	-0.4819	18.3944	2.2993	21.1471
AB	1.9533	6.1790	6.8325	-0.5864	-0.8654	-0.1082	0.0468
C	0.3194	0.4048	-0.1856	9.2202	4.7672	0.5959	1.4204
AC	0.2978	3.9850	-0.2963	9.1742	-0.1714	-0.0214	0.0018
BC	3.2269	0.6192	-0.2629	-0.3842	3.9172	0.4897	0.9591
ABC	2.9522	6.2132	-0.3236	-0.4811	-0.1521	-0.0190	0.0014
D	0.2118	-0.0272	3.6584	2.3246	-0.0456	-0.0057	0.0001
AD	0.1931	-0.1583	5.5618	2.4426	-0.1045	-0.0131	0.0007
BD	2.1146	-0.0216	3.5802	-0.1108	-0.0460	-0.0058	0.0001
ABD	1.8704	-0.2747	5.5940	-0.0607	-0.0969	-0.0121	0.0006
CD	0.3266	-0.0187	-0.1311	1.9034	0.1180	0.0148	0.0009
ACD	0.2926	-0.2442	-0.2531	2.0138	0.0501	0.0063	0.0002
BCD	3.2514	-0.0340	-0.2255	-0.1220	0.1104	0.0138	0.0008
ABCD	3.9618	-0.2896	-0.2556	-0.0301	0.0919	0.0115	0.0005
Total							
22.4902							

Effect Normalization for Case 2⁵⁻¹

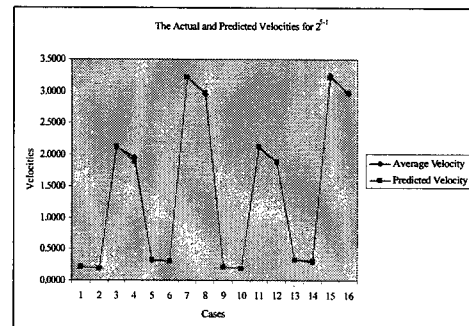
Run	(Sorted)	Norm.
"I"	Effect	Z _i [(I-0.5)/(n-1)]
1	-0.1335	-1.8339
2	-0.1082	-1.2816
3	-0.0214	-0.9674
4	-0.0190	-0.7279
5	-0.0131	-0.5244
6	-0.0121	-0.3407
7	-0.0058	-0.1679
8	-0.0057	0.0000
9	0.0063	0.1679
10	0.0115	0.3407
11	0.0138	0.5244
12	0.0148	0.7279
13	0.4897	0.9674
14	0.5959	1.2816
15	2.2993	1.8339

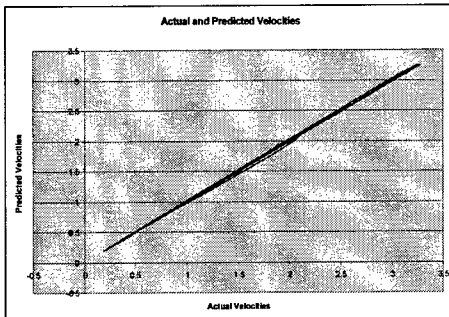
ANOVA Table for Case 2⁵⁻¹

Source	SS	dof	MS	F	p
			=(SS/dof)	=(MS/MSE)	
A	0.0713	1	0.07132919	100.033556	1.5871E-06
B	21.1471	1	21.1471243	29657.172	1.071E-18
C	1.4204	1	1.42038843	1991.98262	7.6684E-13
AB	0.0468	1	0.04680191	65.6359888	1.0541E-05
BC	0.9591	1	0.95905248	1344.99538	5.4044E-12
Error	0.0071	10	0.00071305		
Total	23.6518	15			
R Sq. =	0.99969852				

Predicted velocity = $2.8112745/2 - 0.133537625/2$ A (length) + $2.299300125/2$ B (initial velocity) + $0.59590025/2$ C (diameter ratio) - $0.10816875/2$ AB (interaction of length and velocity) + $0.489656125/2$ BC (interaction of velocity and diameter ratio) + Error.

A	B	C	D	E	Average	Predicted	Error	Perccn.
Length	Init. Vel.	Dia.Ratio	n (Mat.)	m (Fric.)	Velocity	Velocity		Error
-1	-1	-1	-1	1	0.2169	0.2155	0.0014	0.6352
1	-1	-1	-1	-1	0.1897	0.1902	-0.0005	-0.2521
-1	1	-1	-1	-1	2.1117	2.1334	-0.0217	-1.0269
1	1	-1	-1	1	1.9533	1.8917	0.0617	3.1576
-1	-1	1	-1	-1	0.3194	0.3218	-0.0024	-0.7485
1	-1	1	-1	1	0.2978	0.2964	0.0014	0.4628
-1	1	1	-1	1	3.2269	3.2189	0.0080	0.2467
1	1	1	-1	-1	2.9522	2.9772	-0.0251	-0.8488
-1	-1	-1	1	-1	0.2118	0.2135	-0.0018	-1.7929
1	-1	-1	1	1	0.1931	0.1902	0.0029	1.5003
-1	1	-1	1	1	2.1146	2.1334	-0.0188	-0.8871
1	1	-1	1	-1	1.8704	1.8917	-0.0212	-1.1353
-1	-1	1	1	1	0.3266	0.3218	0.0048	1.4725
1	-1	1	1	-1	0.2926	0.2964	-0.0038	-1.2975
-1	1	1	1	-1	3.2514	3.2189	0.0325	0.9991
1	1	1	1	1	2.9618	2.9772	-0.0154	-0.5194





Neural Network “Objective”

- The objective of the second part of the project is to design an artificial neural network based on the method of back-propagation to predict velocities of the extrusion process for the straight converging dies.

- The neural network (NN) is one of the main artificial intelligence (AI) techniques used.
- Artificial neural network is considered as an information-processing system.
- Information processing with neural network consists of analyzing pattern of activity, with learned information stored as weights between node connections.

- A common characteristic is the ability of the system to classify streams of input data without explicit knowledge of rules and to use arbitrary patterns of weights to represent the memory of categories.
- The design of ANN consists of two main stages: TRAINING & TESTING.
- It is trained with numerous examples so as to give reliable results.
- Neural computing offer the advantage of speed once the network has been trained.

Neural Network Results “Training Sets”

	Before DOE Optimization	After DOE Optimization
■ Average Error =	-0.0409	-0.0149
■ MSE (Mean Square Error) =	1.8647	2.2138
■ RMSE (Root MSE) =	1.3655	1.4879
■ Average Absolute Error =	0.8174	0.8562
■ Max. Absolute Error =	10.3798	12.1502
■ Average Percentage Error =	-6.3547	-5.4214

Neural Network Results “Testing Sets”

	Before DOE Optimization	After DOE Optimization
■ Average Error =	-0.0785	-0.0391
■ MSE (Mean Square Error) =	1.6939	2.0887
■ RMSE (Root MSE) =	1.3015	1.4452
■ Average Absolute Error =	0.7942	0.8507
■ Max. Absolute Error =	10.3798	12.1466
■ Average Percentage Error =	-8.7345	-8.0900

- Clear reduction in the error occurred after eliminating the insignificant parameters.

CONCLUSION

- Design of experiments can be used to identify the significance of input parameters, and their overall effect on the model.
- Fractional factorial design can be used to reduce the number of experiments. A more efficient data set can be generated using the information from design of experiments, to train a neural network.

- The neural network model prediction can be used for real-time process control in place of analytical solutions, which can take minutes to hours.
- These advantages are very important in the case of virtual manufacturing.

Incorporating Hybrid Models into a Framework for Design of Multi-Stage Material Processes

Enrique A. Medina*, Daniel A. Allwine**

* Materials & Manufacturing Directorate,
Air Force Research Laboratory,
Wright-Patterson Air Force Base, Ohio

** Austral Engineering and Software Inc., Ohio

Most current applications of software for design of material processes are based on using high fidelity analysis as a substitute for experiments, in a methodology that can still be regarded as trial and error. The state of the art goes somewhat beyond that paradigm by using optimization algorithms to vary the parameters of computationally intensive analysis models to improve the design of individual stages of a manufacturing process. The ongoing work presented here is aimed at creating a tool for preliminary design of products and processes that considers the entire sequence of required processes simultaneously. This is accomplished by viewing the sequence of processes and the product as system that can be optimally designed using formal mathematical techniques. The ultimate goal is to create an easy-to-use software system for integrating material, process, equipment, and cost models, and optimization algorithms into a single environment for preliminary and intermediate design of multi-stage material processes.

Since a large portion of the cost of a system or component is decided early in the design process, the design framework presented will address affordability and sustainability of military and commercial systems by allowing the designer to consider alternative materials and processes and to estimate the influence of design decisions on cost at early design stages. The initial focus of the project is preliminary and intermediate design of sequences of processes used in production of turbine engine components. The proposed design framework is based on a powerful object-oriented, geometry-intensive environment, the Adaptive Modeling Language (AML), and will support different types of models (analytical or numerical) and different mathematical optimization algorithms. Customizable visualization capabilities, web accessibility and distributed, collaborative design are all planned features of the system.

The fundamental difference between this and other simulation-based design tools is the concept that when design is the objective, it is beneficial for the underlying modeling methodology to be formulated taking into account the design function. Basic mathematical analysis enhanced with empirical knowledge can be used to create models that include known design drivers at a level appropriate for preliminary and intermediate design. Optimization algorithms can then be used to vary the parameters of these models in order to solve appropriately formulated design problems that address materials, process, equipment, and cost by means of suitable combinations of objectives and constraints.

The current status of an AML implementation of a three-stage model of a turbine engine disk-manufacturing process will be demonstrated. The model will include material, process, and cost characteristics, and a generalized hill-climbing algorithm will be used to vary model parameters in order to improve a meaningful objective function. Visualization capabilities will include geometry, cost and cost drivers, and measures of optimization algorithm performance.

MPDX™
Multi-Process Design eXecutive:
A Software Tool for Preliminary Design of
Manufacturing Processes

Enrique A. Medina
 Daniel A. Allwine
 Austral Engineering and Software, Inc.

This Small Business Technology Research Project is sponsored by the U.S. Air Force
 Materials Research Laboratory at Wright-Patterson Air Force Base, Ohio, U.S.A.

Motivation: State of the Art in
Manufacturing Design Tools

- Most current applications of software for design of manufacturing processes use high fidelity analysis as a substitute for shop-floor experimentation.
- State of the art uses optimization algorithms to vary the parameters of computationally intensive analysis models to improve individual manufacturing stages.
- This tends to optimize locally and serves to refine an existing design.
- A tool for preliminary design of manufacturing is needed that helps the engineer come up with initial designs. Since cost is committed early, this tool should design for cost as well as properties.

Documented Industrial Needs

Integrated Manufacturing Technology Roadmapping Initiative:
 HTTP://IMTR.ORNL.GOV
 Funded by US DOE, DoD, NSF, DOD (HST)
 Participants: Many U.S. industrial, academic, and technology development organizations

IMTR has identified these recurring themes from industry input that have broad applicability and major impact when solved:

- Modeling and Simulation**
 Multi-scale continuum modeling
 Physics-based models integrated with living knowledge/experience bases
 Models fail-safe, not exception
 Intelligent design & analysis advisors
 MBS as real time enterprise controller
 Smart, self-learning models
 Open, shared repositories & validation centers
 Integrated, robust product and process evolution
 Supporting all domains & applications
 Total, seamless model interoperability
 Real time, interactive, seamless design
- Information Systems for Manufacturing Enterprise**
 Information Systems for Manufacturing Enterprises
 Information-driven seamless enterprises
 Shared knowledge repositories
 Customer / requirements-driven manufacturing
 Mature integrated product / process development
 Totally connected extended enterprises
 Plug-and-play, interoperable systems architectures
 Design and operating systems
 Self-correcting, adaptive control mechanisms
 Self-learning systems
 Integration of mechanical design and other manufacturing design and other manufacturing design
- Manufacturing Processes and Equipment**
 Zero net lay-out waste
 Fast turn-around
 Intelligent control systems
 Innovative breakthrough processes
 Science-based manufacturing
 Intelligent design and process evolution
 Knowledge repositories and validation centers
 Distributed control across extended enterprises
 Engineered materials and surface
 Freeform manufacturing technology

Expected Product Characteristics

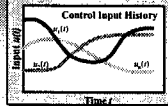
- Tool for preliminary and intermediate process design
- Considers all stages of the manufacturing process
- Rapidly evaluates alternatives for manufacturing and optimizes parameters of selected alternative
- Provides good cost estimation and control at design stages where most of the costs are committed
- Minimizes cost subject to physical constraints given by analytical, reduced models of the physics enhanced with experience-based knowledge
- Incorporates process/business strategies

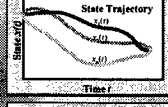
Example of Optimization-Based Design
A common optimal control problem setup

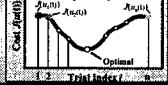
Optimal Control Problem
 Find $u(t)$ to minimize the optimality criterion

$$J(u(t)) = h(x(t_f)) + \int_0^{t_f} g(x(t), u(t), t) dt$$
,
 while satisfying the system state equation
 $\dot{x}(t) = f(x(t), u(t), t), \quad x(0) = x_0$

A Typical Hot Metal Forming Case
 Control inputs, $u(t)$: strain, strain rate, temperature
 System states, $x(t)$: grain size, vol. fraction transformed
 System model, $f(x, u, t)$: microstructural model, deformation model, heat transfer model
 Trajectory criteria, $g(x, u, t)$: regulate microstructure development, minimize deformation heating, stay within workability ranges, integrate equipment characteristics
 Final state criteria, $h(x(t_f))$: achieve final microstructure, achieve final strain

Control Input History


State Trajectory


Optimality Criterion


Example of Optimization-Based Design
Optimal design of an extrusion process
Total strain and final grain size specified

Optimality Criterion

$$J(\tilde{e}(t)) = 10(\tilde{e}(t_f) - 2.0)^2 + \int_0^{t_f} (d(t) - d_{desired})^2 dt$$

System State Equation

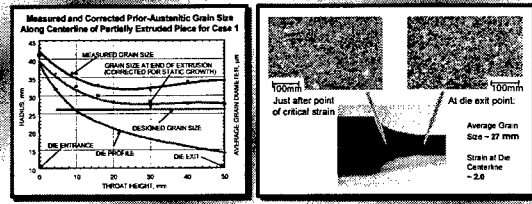
$$\begin{bmatrix} \dot{d} \\ \dot{\tilde{e}} \end{bmatrix} = \begin{bmatrix} f_1(T, u) \\ f_2(x, \tilde{e}, T, u) \\ f_3(\tilde{e}, T, u) \end{bmatrix}$$

EXTRUSION OF PLAIN CARBON STEEL RODS

CASE	DESIRABLE GRAIN SIZE (mm)	RAM VELOCITY (mm/s)
OPTIMAL SOLUTION 1	0.025	8.43
OPTIMAL SOLUTION 2	0.025	10.0
OPTIMAL SOLUTION 3	0.025	15.0
CONSTRAIN TO MAXIMIZE	0.025	10.0

Example of Optimization-Based Design of Material Processes

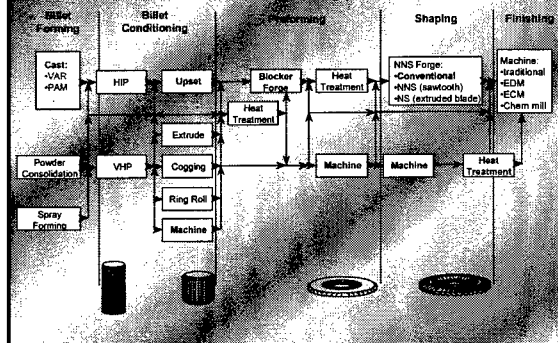
- Die profile and extrusion speed calculated through optimization
- Experimental results agreed with design predictions based on simple, dynamic models



Interactive Simulation System for Design of Multi-Step Material Processes: SBIR Project

- Develop framework for integration of models for materials, processes, equipment, geometry and cost, and design and optimization algorithms.
- Develop process design models in a model library approach. This will allow the creation of a library of model features that can be used to construct new models.
- Develop interactive visual user interface for design system.
- Demonstrate and validate effectiveness on Air Force problem (focus is manufacturing of turbine engine disks).

Typical Alternatives for Manufacturing of Ti-64 Turbine Engine Disks



Modeling for Properties and Cost

- Models are required for estimating
 - part properties
 - manufacturing cost
- Modeling efforts
 - Materials
 - geometry - parametric
 - field variables - global or local basis functions, discrete points
 - Process
 - Equipment
 - Cost
- Models must relate physics to properties and cost
- Models should execute rapidly allowing evaluation of many alternative designs
- Use object-oriented methodologies

First Process Models Example

- Forging model for a disk with cross section given by



- Geometry is parameterized (radii and heights)
- An enhanced slab-analysis method is used to parameterize
 - for each of the four slabs, an average value for
 - accumulated strain
 - strain rate
 - temperature
 - pressure and load
- Process parameters can be changed
- Flat die, open die, and closed die forging cases modeled
- Multi-stage forging modeled (intermediate geometries accepted)

First Cost Model: Brief Overview

Assumptions

isothermal forging of Ti-64 for rotors on a particular press.
Initial beta-quench operation with alpha-beta forging desired.
Fixed number of rotors produced; dedicated production equipment.

Cost distribution

1. Material
2. Upset Forging
3. Blocker Forge
4. Close Die Forging
5. Extrusion
6. Preform Machining
7. Rough Machining
8. Finish Machining
9. Heat Treatment
10. Ultrasonic Inspection

Sample rules

In addition to normal rules based on labor, energy, and material:
Probability of cracking is based on stable processing window.
Die cost uses approximate life of a T2M die.
Heat treatment cost varies with volume fraction globularized of pre-heat treatment part.

Some shortcomings

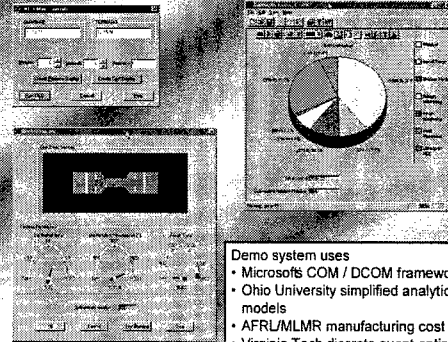
Level of detail should be made uniform.
Quantities such as the production run should be made variable.
Cost should be modeled on a per-object basis.
Not all cost drivers have been considered.
Should consider equipment as a variable.

Phase I Proof of Concept: Problem

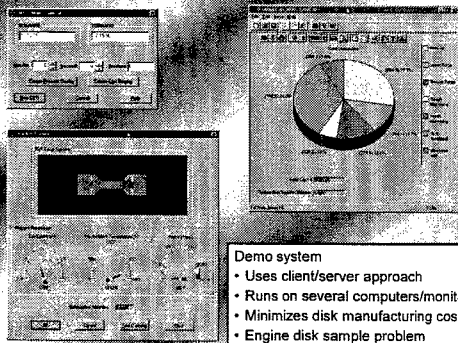
- Alternative routes for manufacturing a sample disk
- Workpiece and die geometries are parameterized
- Optimization algorithm varies workpiece, die, and process parameters to minimize cost
- Choice of sequence was not an optimization variable



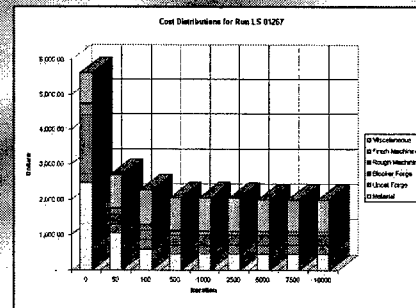
Phase I Proof of Concept



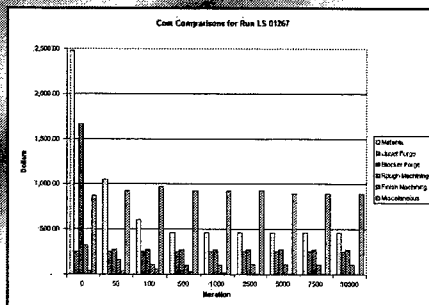
Phase I Proof of Concept



Optimization Results for a Run of a Sample Problem



Optimization Results for a Run of a Sample Problem



The Adaptive Modeling Language

- For a mature design product, a solid foundation is required
- TechnoSoft's AML has been selected for building the integration framework
- AML is object oriented and feature based
- AML provides a solid integration framework
- MEDXTM is using AML to integrate
 - simplified analytical models and knowledge-based models
 - optimization algorithms
 - design rules
 - custom visualization mechanisms
 - cost models
- into a system for preliminary design of multi-stage manufacturing processes

AML Capabilities

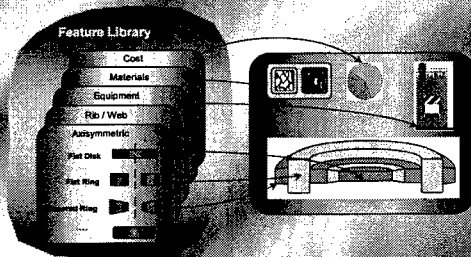
- Unified product model
- Full support of UNIX and Windows NT
- Single Underlying Object-Oriented Architecture
- Open Architecture for Foreign Applications Seamless Integration
- Common Syntax throughout the different Modules
- Real Time Dependency Tracking
- Demand Driven Computation
- Full Support of IGES (STEP)
- Support of various geometric modelers with full model compatibility

AML Distributed-Computing Capabilities

- Network Distributed Architecture
 - OS Independent
- CORBA Compliant
- STEP Compliant
- Environment supporting real time
 - Interactive modeling of complex systems
 - Dynamic linking of the various engineering processes
 - Distributed models
 - Remote access via standard browsers

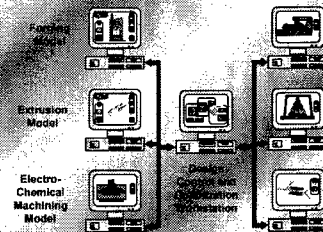
MPDX™ Facilitated Model Building

- Models are built from libraries that contain both empirical knowledge and physics



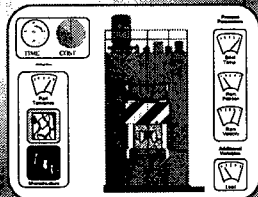
MPDX™ Automated, Distributed Design Capabilities

- The communication medium can be the operating system, an intranet, the Internet, or the Web



MPDX™ Interactive Design Capabilities

- Change parameters and quickly evaluate design
- Get real time information about impact of design decisions (material, process, equipment) on manufacturing cost and product properties
- Run optimization to find best processing sequences and parameters



Summary

- Tool for preliminary design of multi-stage manufacturing processes
- When design is the objective, the models should be appropriate for design.
- Basic mathematical analysis and empirical knowledge are used for modeling that includes dynamic effects.
- Optimization is used to find best processing sequence and parameters
- State-of-the-art integration framework with distributed computing capabilities

Hybrid Modeling for the Interdisciplinary Design of More Affordable Systems

James W. Poindexter, Gerald R. Shumaker and Brian A. Stucke

Materials and Manufacturing Directorate,
Air Force Research Laboratory, Wright-Patterson AFB, OH, USA

This paper will address current results and future requirements for system modeling and simulation with an emphasis on enhancing physics-based shop floor models while developing a common data model to support the integration of shape, material, process and performance characteristics. A method of "interdigitation", which tightly couples an engineers' domain specific knowledge with domain independent search techniques such as numerical optimization, genetic algorithms and simulated annealing, leverages the strengths of each technique to provide a more holistic approach to system design.

Developing and building systems in the 21st century will require an extensive effort in modeling and simulation to meet the cost and performance requirements for new Department of Defense (DoD) systems. Advanced, multi-role systems like the Joint Strike Fighter (JSF) are striving to achieve strict cost objectives in the face of shifting system performance requirements. To balance these cost and performance goals, increased emphasis is being placed on the use of emerging technologies such as multi-disciplinary modeling and simulation and a collaborative optimization environment to accelerate and enhance cost decisions in the design of an aircraft.

New initiatives involving "cost-as-an-independent-variable" (CAIV) are dramatically changing the Air Force and DoD acquisition process from performance parameters and pure acquisition cost to assessing a system's total life cycle and total ownership cost. CAIV will require extensive communication between engineering and manufacturing disciplines to understand the complexity of cost impacts on a system which is greatly enhanced through the use of modeling and simulation of both the system and its environment from manufacture through to retirement. To understand the relationship between performance and manufacturability requires detailed knowledge of the entire scope of processes used to develop and support a system from the macro to micro scale.

Several DoD research and development initiatives are focusing on enhancing the functional relationship between critical design, engineering data and the tools that model the manufacturing fitness of a component. One such initiative is the Simulation Assessment and Validation Environment (SAVE), an open system architecture built on the Common Object Request Broker Architecture (CORBA), which enables seamless sharing of data across a wide array of commercially-available modeling and simulation tools in support of system design. The impact of such an integrated suite of simulation tools applied to the JSF, is estimated to save \$3 Billion in life cycle costs.

The collaborative optimization environment afforded by SAVE offers a seamless development environment necessary for engineers to model, analyze, and optimize complex products and processes that produce them. This environment seeks to integrate design performance tools such as finite element analysis and computational fluid dynamic models with cost models associated with particular materials and processes. As a design configuration evolves, the multi-disciplinary team gets immediate feedback from a manufacturing process database on how well the design fits within manufacturing capabilities, and whether the performance of the component is optimized. Ultimately this type of collaborative environment will interface with a DoD battlefield modeling and simulation environment to provide the acquisition community (government and contractors) with the ability to represent more fully a product's key characteristics and performance capabilities with higher fidelity.



The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

Hybrid Modeling for the Interdisciplinary Design of More Affordable Systems

Brian A. Stucke

AFRL/MLMS
 2977 P Street Suite 6
 Wright Patterson AFB, OH 45433-7739
 tel: (937) 255-4623
 fax: (937) 656-4269
 brian.stucke@afrl.af.mil

James W. Poindexter

AFRL/MLMS
 2977 P Street Suite 6
 Wright Patterson AFB, OH 45433-7739
 tel: (937) 255-7371
 fax: (937) 656-4269
 james.poindexter@afrl.af.mil

Briefing Period: 01/02/1999

Hybrid Modeling for the Interdisciplinary Design of More Affordable Systems™ - SPANWEE Brief - Slide 1

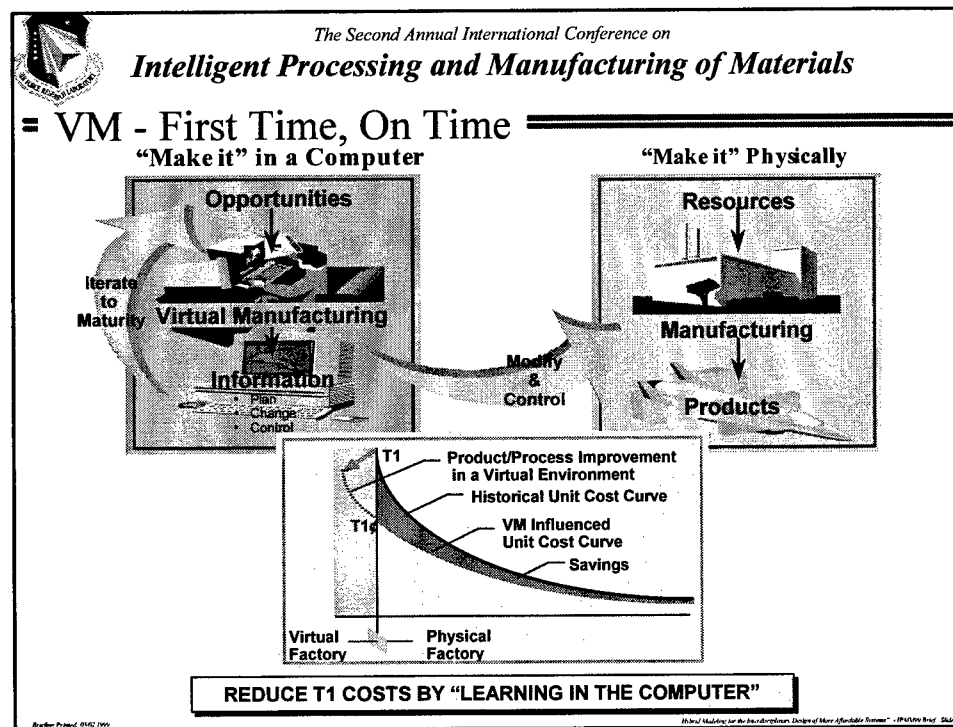
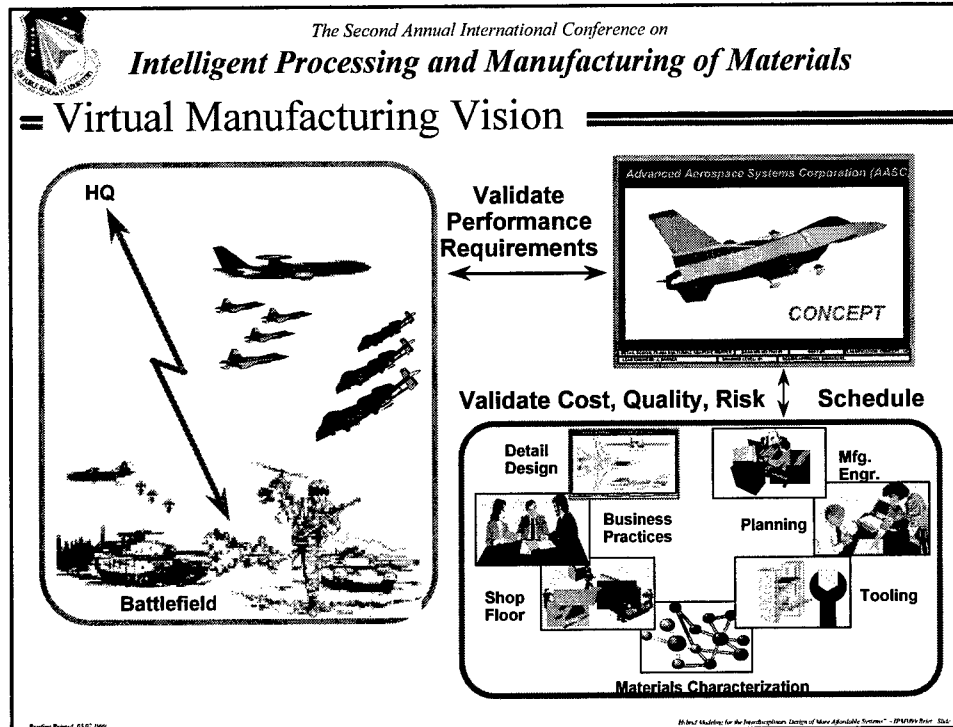


The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

Overview

- Introduction
- Virtual Manufacturing Vision
- Multi-discipline Modeling & Simulation
- Collaborative Engineering Environment
- Simulation Assessment Validation Environment
- Benefits and Conclusions
 - Integrated Manufacturing Simulation for Affordability

Briefing Period





The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

= Manufacturing Modeling & Simulation Challenges =

- Information Representation
 - “Continuum Modeling” - micro to macro level
 - Legacy Models
- Model/Data Integration
 - Multiple view Product Data
 - Smart Product Models
- Business Practices & Operations
 - Implementing M&S Technologies into the Product Realization Process
 - Business Case
 - Benefits/Metrics
 - M&S used for Real Time Enterprise Control
 - Links to Simulation Based Acquisition Community
- AFRL/MLMS Research Projects addressing these Challenges
 - Collaborative Engineering Environment - CEE
 - Simulation Assessment Validation Environment - SAVE
 - Integrated Manufacturing Simulation for Affordability - IMSA

Working Project: 0102 100

Model Modeling for the Interdisciplinary Design of New Affordable Systems™ - IPIMMS Model Slide 1



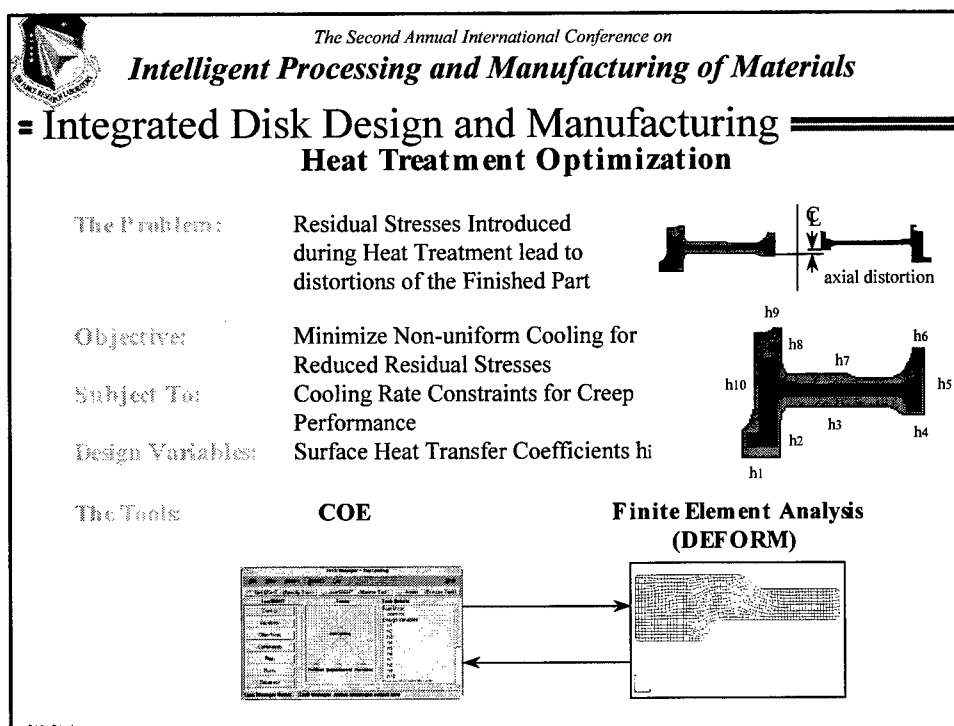
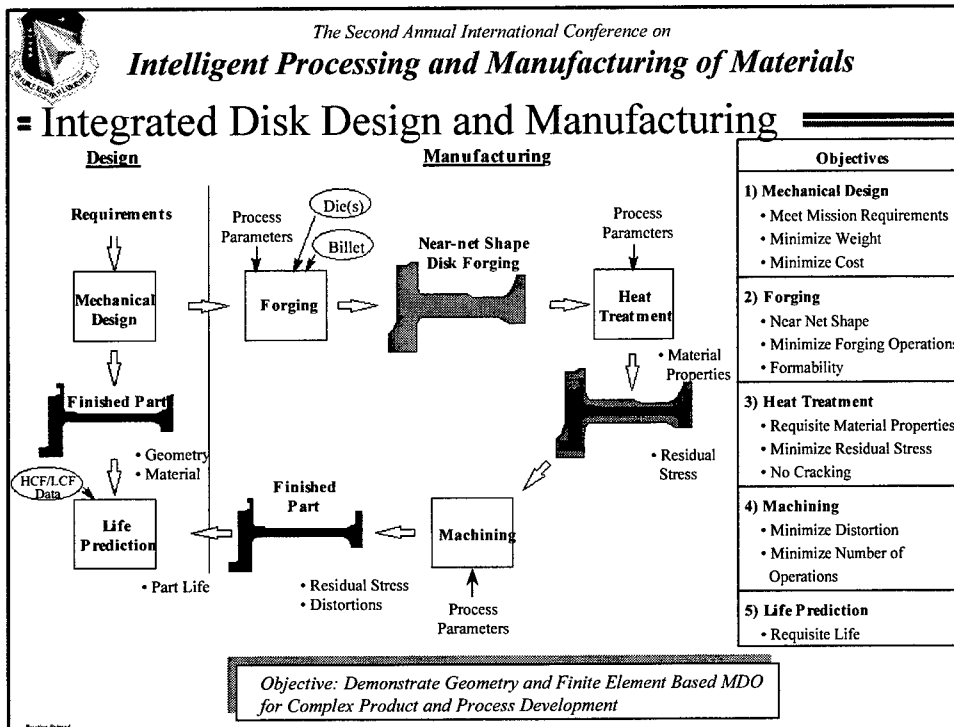
The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

Collaborative Engineering Environment CEE

Working Project: 0102 100

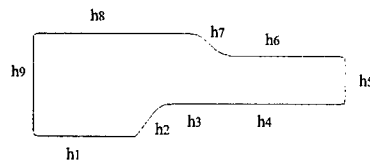
Model Modeling for the Interdisciplinary Design of New Affordable Systems™ - IPIMMS Model Slide 2



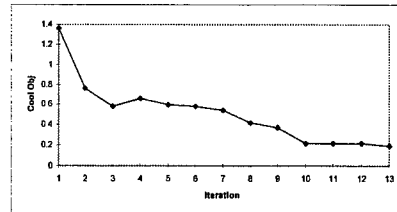


The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

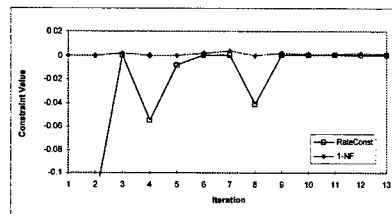
= Integrated Disk Design and Manufacturing
Optimization Results for Generic HPT Disk



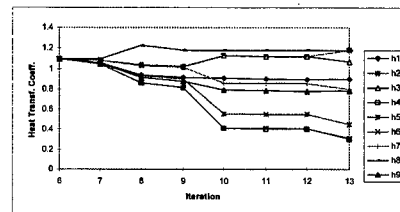
Design Variables



Objective Function History



Constraint History



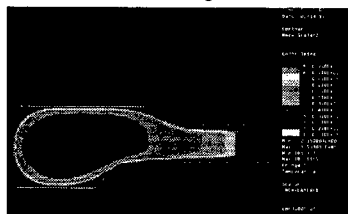
Design Variable History $[10^{-4} \text{ BTU} / (\text{in}^2 \text{ s } ^\circ \text{F})]$



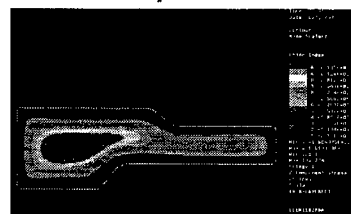
The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

= Integrated Disk Design and Manufacturing
Cooling Rate and Stress Results

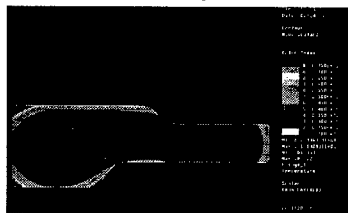
Initial Cooling Rates



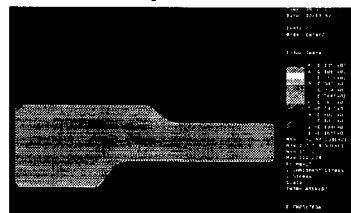
Initial Hoop Stress Distribution



Final Cooling Rates



Final Hoop Stress Distribution



- Cooling Rate Distribution more Uniform
- Cooling Rate Target Met

Maximum Stresses Reduced by about 80 %



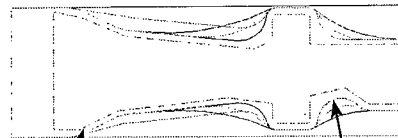
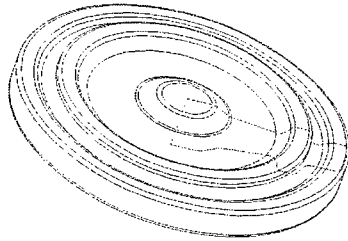
The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

= Integrated Disk Design and Manufacturing

Disk Near-Net-Shape Forging Optimization

Representative Shapes



Sonic Shape Current Pancake Shape Most Aggressive Near-Net-Shape

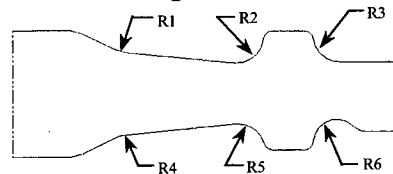
Objective

Determine minimum weight for the near-net-shape forging of turbine disk

Constraints

- max press load \leq allowable value
- max strain rate \leq allowable value
- min corner and fillet radius
- sonic shape constraint

Design Parameters



The Second Annual International Conference on

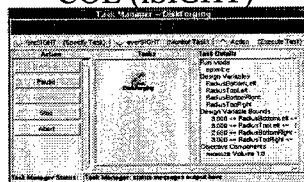
Intelligent Processing and Manufacturing of Materials

= Integrated Disk Design and Manufacturing

The Approach To Forging Optimization

COE (iSIGHT)

- Select design variables
- Display geometry
- Update models
- Compute volumes
- Export geometry

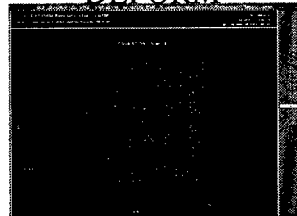


- Generate mesh
- Generate BC
- Start DEFORM run
- Monitor DEFORM run
- Postprocess results

Unigraphics



DEFORM



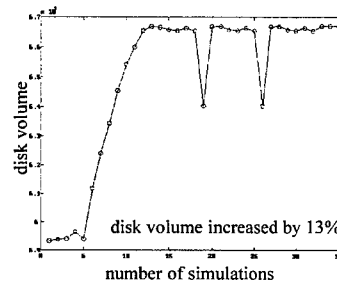
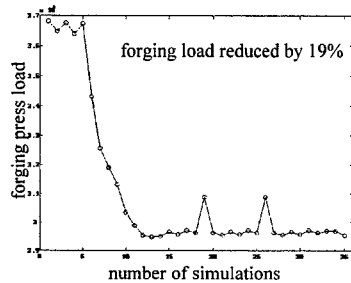
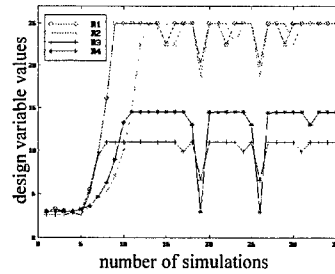
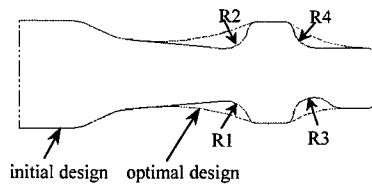


The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

= Integrated Disk Design and Manufacturing Optimization Results for Single-Step Isothermal Forging

Disk Design Parameters



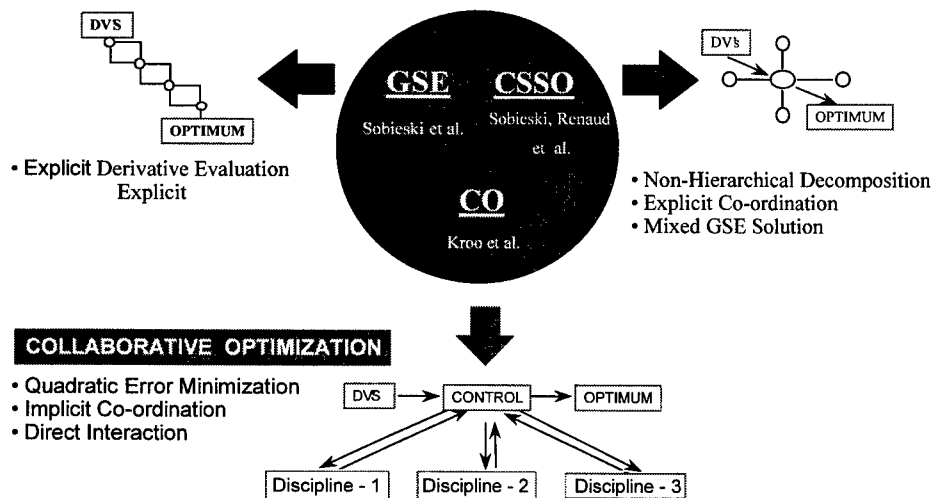
The Second Annual International Conference on

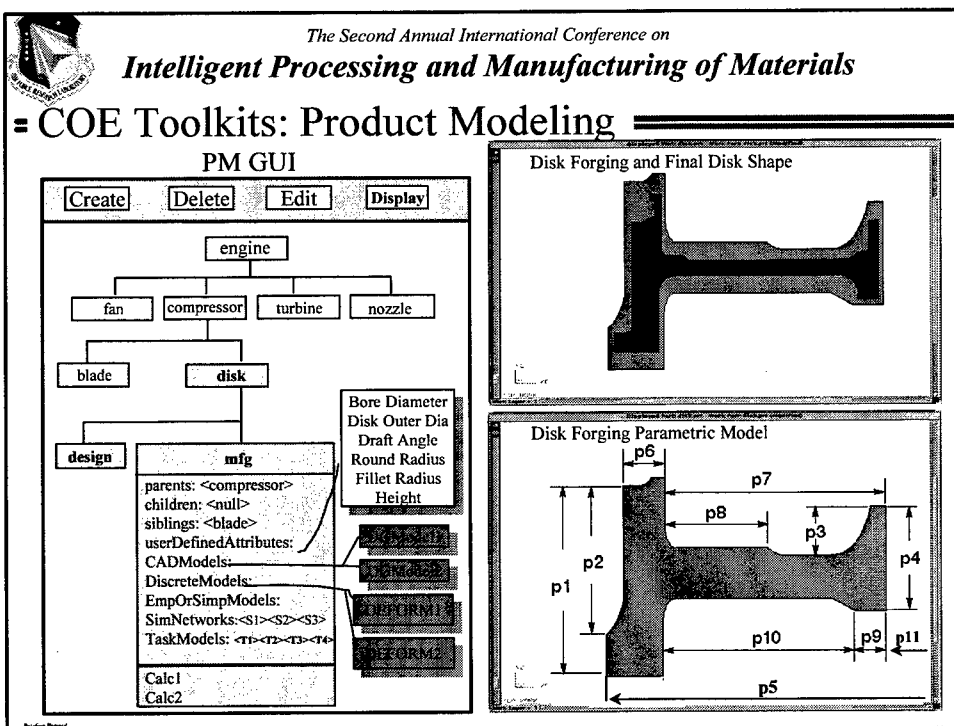
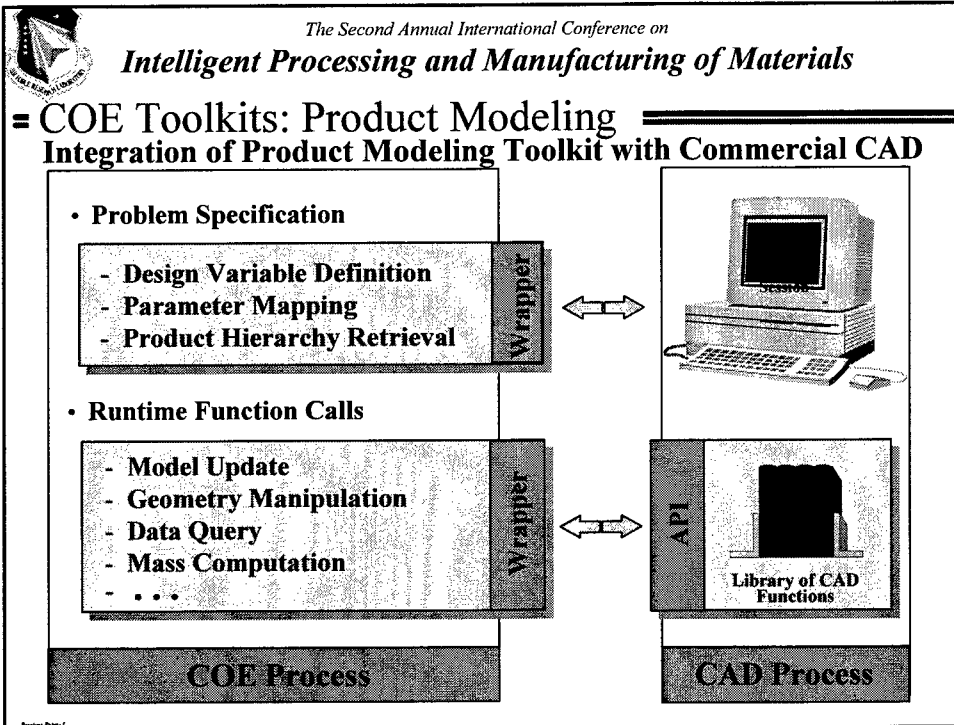
Intelligent Processing and Manufacturing of Materials

= COE Toolkits: Optimization Multidisciplinary Design Optimization Methods

GLOBAL SENSITIVITY EQUATIONS

CONCURRENT SUBSPACE OPTIMIZATION







The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

Simulation Assessment Validation Environment SAVE

Working Paper 01021999

Virtual Modeling for the Investigation: Design of More Affordable Systems - JPMODs Rev. Slide 17



The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

= Problem Statement

- **Accurate & Timely Knowledge Of Manufacturing Impacts Critical For:**
 - * process selection and planning
 - * tooling/fixtures design & verification
 - * assembly planning & control
 - * cost of performance analysis
- **There is currently no toolset available that*:**
 - * bridges the gap between different manufacturing modeling & simulation tools
 - * provides comprehensive visibility into manufacturing for design supporting design decision impact analysis
 - * allows accurate representation of individual processes and their interaction to support realistic and timely trade-off analysis

• Reducing Costs for JSF Requires Tools that Accurately Represent and Use Manufacturing Information During Design

* From User & Technical Workshops (over 140 participants)

Working Paper

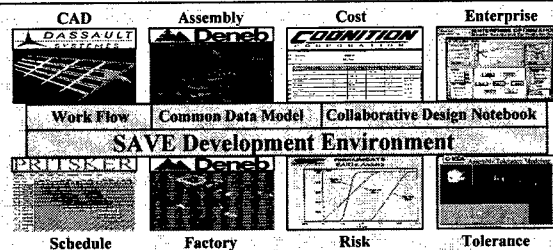


The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

= SAVE Program Objective

Demonstrate, validate and implement integrated modeling and simulation tools and methods to assess the impacts on manufacturing of product/process decisions supporting low risk, affordable transition of weapon systems technology from design to EMD



- Assess Mfg Impact of Product/Process Design Decisions
- Supports Assessment of Cost as an Independent Variable
- Validates Product & Process Prior To Design Release
- Reduces Risk In Transition To Production
- Estimated 2-3% Life Cycle Cost Reduction (\$3B on JSF)



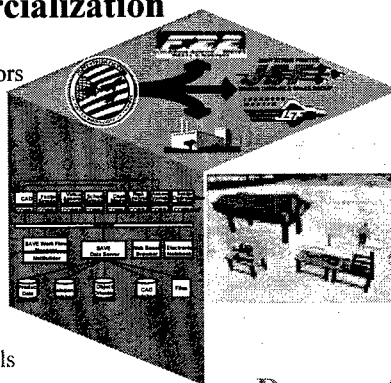
The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

= SAVE Focus Areas

Implementation/Commercialization

- WSC Beta Testing
- SAVE availability to other JSF contractors
- Commercial end product embraced by participating SAVE tool vendors



Development



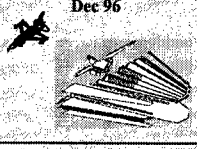


- Plug and play via tool wrapping
- Open architecture implementation
- Integration of commercial M&S tools

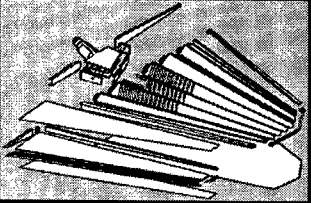
Demonstration

- Identification/Assessment of appropriate metrics
- Early introduction of meaningful manufacturing information
- Application of VM technologies into product development process

The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

= SAVE Phase I Demonstration

 IP/PT Team  Integrated Tool Set	INITIAL DEMO Dec 96  Upgrade / Mod Scenario	INTERIM DEMO July 98  Design / Mfg Trade Study Scenario	FINAL DEMO June 99  Assy Optimization Scenario	Videos & Industry Review
---	---	---	---	-------------------------------------


**Initial Demonstration
F-16 Horizontal Stabilizer**

**Optimize Implementation of F-16
Horizontal Stabilizer Modification**

- **Recommendation:** Net Trim Skin on all Five Sides and Eliminate FOUR Assembly/Disassembly Steps per Skin.

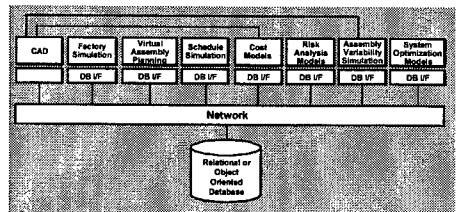
Est Savings: \$114K

Boeing Patent: 5,142,799 Federal Modeling for the Interdisciplinary Design of Man-Affordable Systems™ - IP/PT/PPR Brief Slide 21

The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

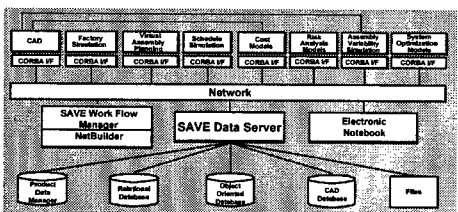
= SAVE Phase II Development Objective

Refine the SAVE Development Environment using a Common Object Request Broker Architecture (CORBA) Based Approach to Simplify Tool Integration and Reduce Implementation Risks.



Phase I Approach

- SAVE Proof-of-Concept
- Data Integration via Common Database
- Required User Maintenance of SAVE Specific Data Base



Phase II Approach

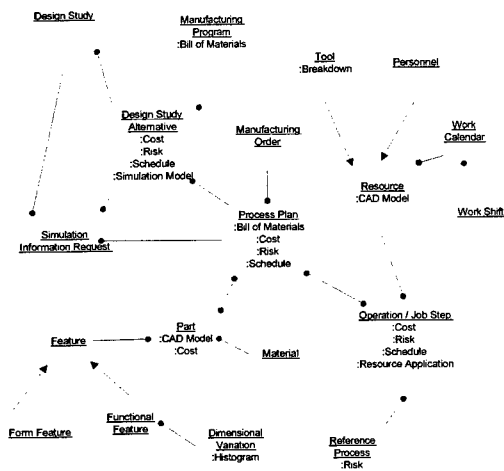
- Leverages Latest Commercial OO Technologies
- Simplifies Plug-n-Play
- Promotes Access to Existing User Legacy Databases
- Reduces Data Ownership & Maintenance Costs

Boeing Patent: 5,142,799 Federal Modeling for the Interdisciplinary Design of Man-Affordable Systems™ - IP/PT/PPR Brief Slide 22

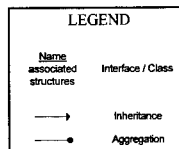


The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

= SAVE Data Model



- Use Numerous Sources of Information
 - SAVE Tool Input/Output Specification Release 1.0
 - SAVE Tool Vendors
 - SAVE Tool Users
 - Lockheed Martin Manufacturing Engineers
- Derive from Top Down Bottoms Up and View
- Capture Shared Data for Classes of Simulations in SAVE Environment
- Capture Data Necessary to Assess the Outcome of a Design Trade Study
- Review Model with Team and Outside Organizations



Brother Project 8/12/1999

Model Modeling for the Remanufacturing Design of More Sustainable Systems™ - (IPIMMS) Brief Slide 21



The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

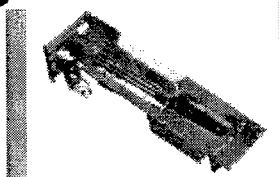
= SAVE Phase II Interim Demonstration

Goals

- Show toolset/database integration
- Demonstrate F-22 IPT Usage
- Show savings from using tools

Success Criteria

- Plug and Play Capability
- Multiple Data Source Access
- # of IPT's Using SAVE/Total # of IPT's
- # of Tools Used
- # of Decisions Affected
- Design to Cost Data Accuracy
- Mfg Lead Time Reduction
- Design Change Reduction
- Scrap, Rework & Repair Reduction
- Inventory Turn Increase
- Fab & Assembly Inspection Reduction



Interim Demonstration

**Perform Design
Manufacturing Trade
Studies**

F-22 Gun Port Assembly

Brother Project 8/12/1999

Model Modeling for the Remanufacturing Design of More Sustainable Systems™ - (IPIMMS) Brief Slide 24



The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

= SAVE Phase II Demo Problem Description

F-22 Gun Port Redesign Activity

- Blast from gun is eroding forward skin and surrounding structure
- Team identified three options
 - Metallic Skin
 - Split Skin (Composite and Metallic)
 - Replaceable Insert and Cover
- Insert and cover option selected for performance reasons
- SAVE evaluated candidate using the virtual manufacturing environment

Reading Panel

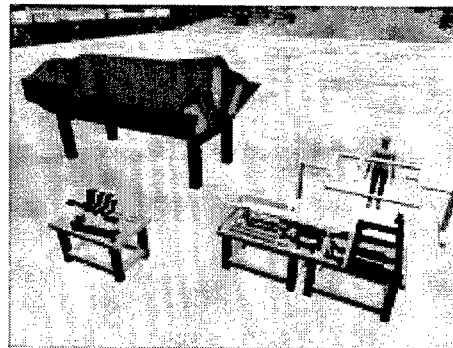


The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

= SAVE Phase II Demo Trade Study Results

- Bottleneck identified in factory
 - Overtime or additional shifts
 - Additional mate tools
- Original process plan not practical
 - Ergonomic issues
 - Assembly sequence
- Revised process plan
 - Additional mate tools
 - Modified assembly sequence
- Titanium vs. Stainless
 - Cost and risk equivalent for both
 - Make selection based on performance criteria



Reading Panel



The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

= SAVE JSF LCC Benefits

PRODUCT/PROCESS METRIC	SAVE IMPACT TO METRIC (%)		UNIT COST SAVINGS (\$K)	
	F-22	JSF	F-22	JSF
Design to Cost Data Accuracy	25	12	109	237
Lead Time Reduction	N/A	10	N/A	155
Design Change Reduction	15	28	326	325
Scrap Rework & Repair Reduction	15	11	80	200
Process Capability	10	5	97	107
Inventory Turn Increase	5	2	25	23
Fab & Assy Inspection Reduction	13	6	79	85
TOTAL			716	1,100

Based on Full Implementation of SAVE Throughout Weapon System

Boeing Patent #712,199

Model Modeling for the Development Design of More Affordable Systems - IPANW/Boeing Slide 27



The Second Annual International Conference on

Intelligent Processing and Manufacturing of Materials

= Integrated Manufacturing Simulation for Affordability



Contracting Strategy - BAA

Business Strategy - Major initiative with multiple tasks including: development, demonstration implementation and marketing.

Start Date: FY 00

Project Engineer: James Poindexter

	FY	00	01	02	03	04	Tot
Gospel		0.5	0.5	4.0	3.0	3.0	11.0
Requested		0.5	2.0	4.0	3.0	1.5	11.0
Cost Share		0.5	1.0	3.0	2.0	1.0	7.5

OBJECTIVE

- Increase Producibility/Affordability considerations early in development
- Improve upstream prod/proc decisions to reduce mfg risk, cost & time

APPROACH

- Focus on Design/Manufacturing Links earlier in the product life cycle
- Tool Development & Product Data Model Representation
- Validation Through Multiple Demonstrations
- Extensive Commercialization & Tech Trans Program

DELIVERABLES

- Concept Development Software Tools and Product Models
- W/S Focused Demonstrations & Videos
- Commercial Level Software & Implementation Plans
- Benchmarks for the Standards community

- **Customers** - AF W/S systems integrators, their electrical and mechanical subsystem suppliers & their s/w tool vendors

Benefits -

- 60% cost & time reduction in physical prototypes
- 90% reduction in errors during dev. & production
- 50% reduction in support costs

Implementation -

- Major contractor/supplier demonstrations
- Commercializable Vendor Products/Services

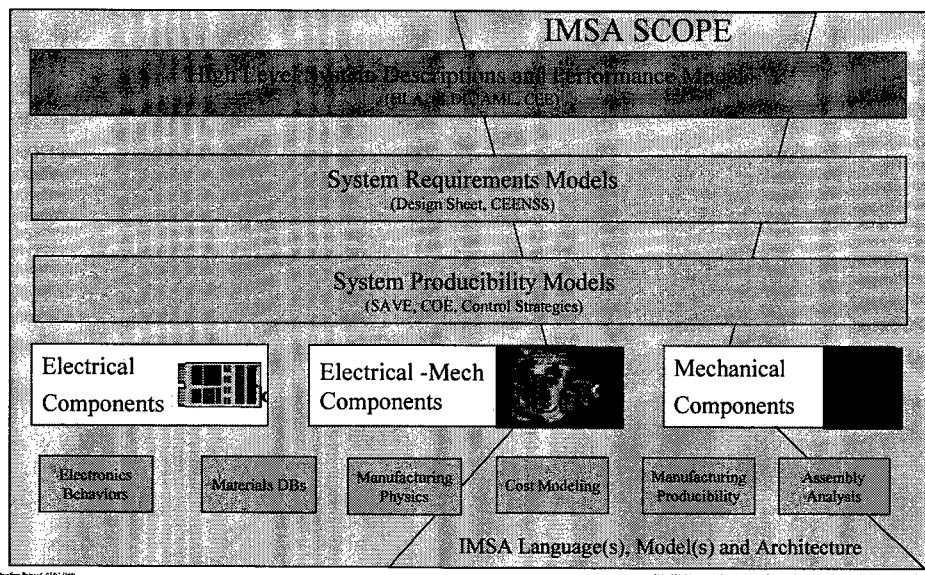
- **Related Efforts** - SAVE, CEENSS, DARPA initiatives, (See Roadmap)

Boeing Patent



The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

■ IMSA Approach & Scope



The Second Annual International Conference on
Intelligent Processing and Manufacturing of Materials

■ IMSA Benefits

- Multi-Disciplinary approach to designing/manufacturing and fielding a system which meets the requirements at desired LCC
- Increase accuracy and number of manufacturing trades on weapon system or components - - 10X, 5 years to 6 months
- Capability to look at critical key characteristics of system and assess manufacturability, cost, risk and schedule
- Depth and Fidelity of data to enhance look forward analysis of trade decision implications over the life of the system - - Reduce costly errors - 100X cost avoidance

Revised Printed: 03/02/1999

Revised Printed: 03/02/1999

NASA Johnson Space Center
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

Hybrid Modeling for Testing Intelligent Software for Lunar-Mars Closed Life Support

Jane T. Malin
NASA Johnson Space Center

2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 1

NASA Johnson Space Center
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

Overview

- Hybrid models for evaluation of intelligent control systems for processing plants
- Four types of models
- Integrating model types in CONFIG modeling and simulation tool
- Test of intelligent software for control of product gas transfer in Lunar-Mars Life Support Test

2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 2

NASA Johnson Space Center
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

Discrete and continuous control in intelligent management of processing plants

- Manage a system where processors convert resources to products
 - Store and transport resources and products
 - With plans and schedules
 - Configure the system and its component modes to transfer, process and store the resources and products
 - With sequences and procedures
 - Control processors that produce products from resources
 - With discrete or continuous control
 - Manage instrumentation and control subsystems
 - Configure and calibrate with manual or discrete control

2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 3

NASA Johnson Space Center
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

Example: Product Gas Transfer in Phase III Lunar Mars Life Support Test

- 90-day closed life support test at NASA Johnson Space Center in fall of 1997
- Four crew members in 20 foot diameter chamber, staged wheat crops in plant growth chamber
- Air and water recycling, using physico-chemical and biological processors
- Intelligent autonomy software for control of storage and transfer of oxygen and carbon dioxide

2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 4

NASA Johnson Space Center
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

Product Gas Transfer in the Phase III Test

The diagram illustrates the product gas transfer system. On the left, the '20 Foot Chamber' contains a 'Sleeping Quarters', 'Air Revitalization System', 'CO₂ Accumulator', 'Water Recovery System', 'Kitchen', and 'Exercise Room'. On the right, the 'Variable Pressure Growth Chamber' is connected to an 'Incinerator'. Gas flow is indicated by arrows: O₂ flows from the 20 Foot Chamber to the O₂ Storage tank, then to the O₂ Concentrator, and finally to the Incinerator. CO₂ flows from the Incinerator back to the 20 Foot Chamber. An 'AIR LOCK' is shown between the two chambers.

2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 5

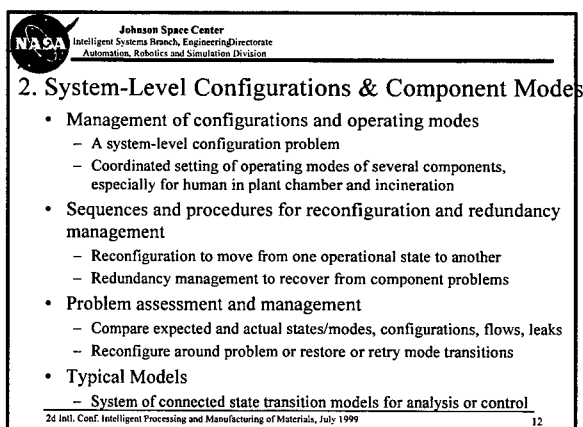
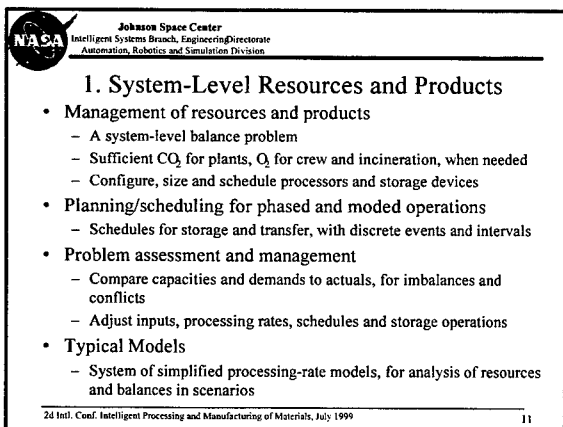
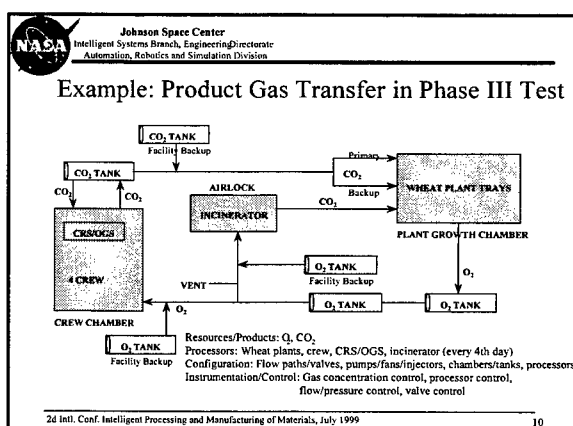
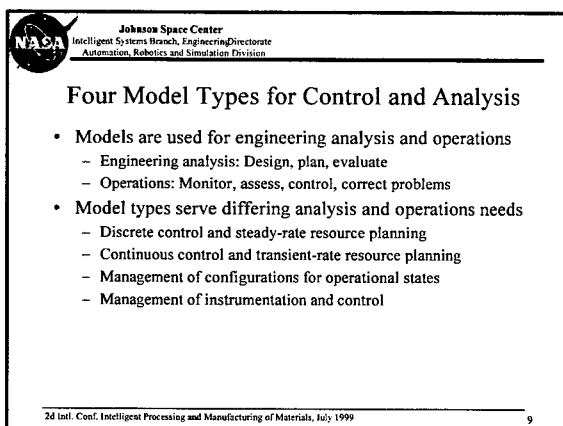
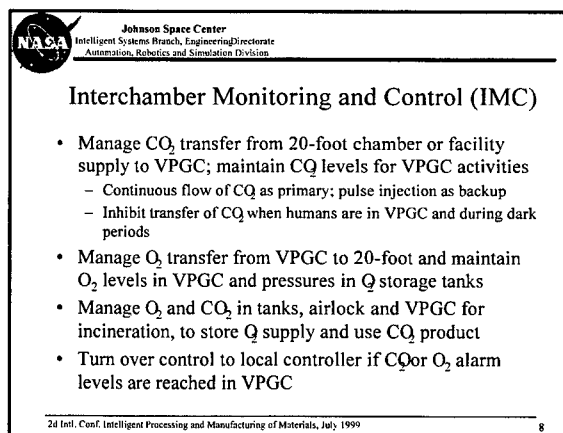
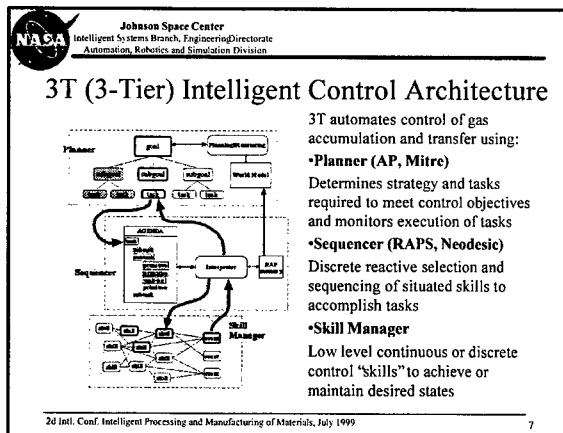
NASA Johnson Space Center
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division


Gas Transfers in Lunar Mars Life Support Test

Complex Automated Control Problem

- Control product gas transfer throughout a 90-day manned test with 4 crew in the 20-foot chamber and staged wheat crops in the plant chamber
- Accumulate product gases (oxygen and carbon dioxide) and transfer them among the 20-foot crew habitat, plant chamber, and incinerator in a closed recycling system
- Adaptively reconfigure and transfer gas among multiple reservoirs in response to predicted needs, observed usage and problems with elements of the system

2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 6




 **Johnson Space Center**
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

3. Processes and Processor Performance

- Management of processor performance
 - Adjust processor parameters to achieve desired performance
 - Wheat plant performance was affected by several variables
- Continuous and discrete control of processors
 - Continuous control in nominal operations and for degradations
 - Discrete control for safety and problem management
- Problem assessment and management
 - Continuous state estimation, comparison to setpoints, assess operating conditions
 - Adjust parameters or operating conditions, reconfigure
- Typical Models
 - Differential equation or algebraic models, analytic or empirical


2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 13

 **Johnson Space Center**
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

4. Subsystems for Instrumentation and Control

- Management of instrumentation and control subsystems
 - Instrumentation and control subsystems are systems themselves
 - Multiple controllers throughout the product gas transfer system
- Schedules, sequences and procedures and parameters for control regimes
 - Activate, deactivate and switch control regimes
 - Adjust setpoints and estimators, calibrate
- Problem assessment and management
 - Assess drifts, mode problems, bad inputs, operating conditions
 - Correct inputs or operating conditions, retry, reconfigure, adjust models, estimators
- Typical Models
 - State transition models for modes of control or control regimes


2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 14

 **Johnson Space Center**
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

Summary of model types for system management

- Management of resources and products
 - System of simplified processing-rate models
- Management of configurations and operating modes
 - System of connected state transition models
- Management of processor performance
 - Differential equation or algebraic models, analytic or empirical
- Management of instrumentation and control subsystems
 - State transition models for modes of control or control regimes


2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 15

 **Johnson Space Center**
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

Integrating model types in CONFIG

- Modeling Product Gas Transfer for the Phase III Lunar Mars Life Support Test
 - Model was developed to perform simulation-based testing and validation of the autonomy software
 - System with all four types of elements required all four model types
- CONFIG is an object-oriented enhanced discrete event simulation system
 - Component models have discrete modes, with continuous dynamic performance within modes
 - Activity models have discrete phases for control regimes, with discrete or continuous control within phases
 - Components are connected in reconfigurable paths of flow and effort


2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 16

 **Johnson Space Center**
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

CONFIG Discrete Event Simulation

- Purpose: Fault-injectable simulations for evaluation of autonomy software
 - Simulate hardware and biological systems and control for interactive dynamic testing and evaluation
 - Provide integrated system-level testing in multiple long-duration scenarios, faster than real time
- Purpose: Analyze effects of events on a system
 - Local and remote, behavioral and functional effects
 - Time course of effects
 - How and where effects are detectable
- Both qualitative and quantitative modeling within discrete event simulation framework

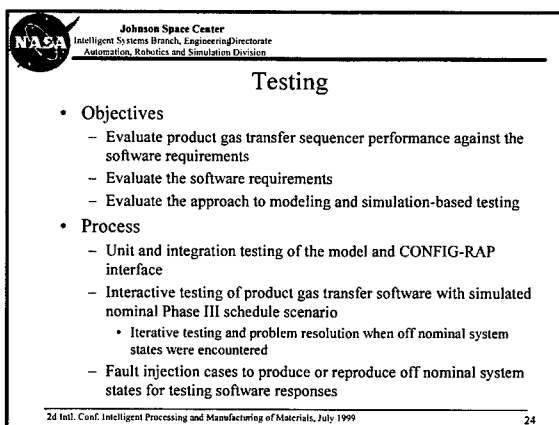
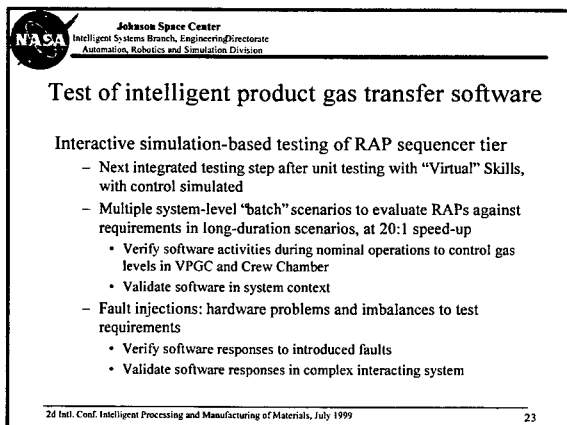
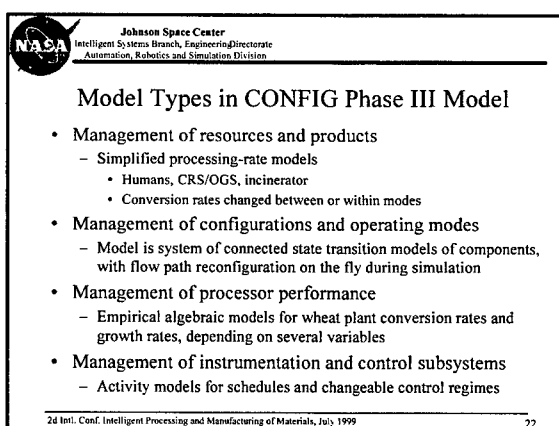
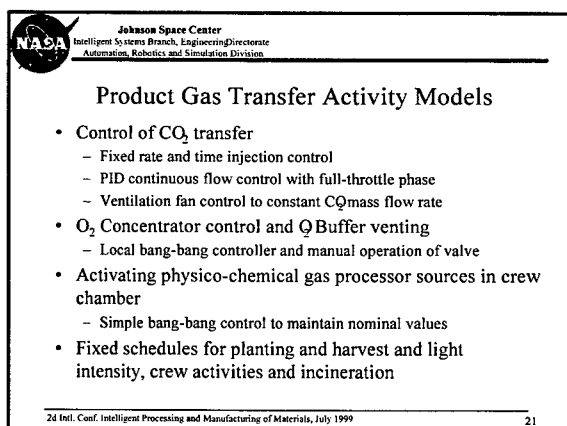
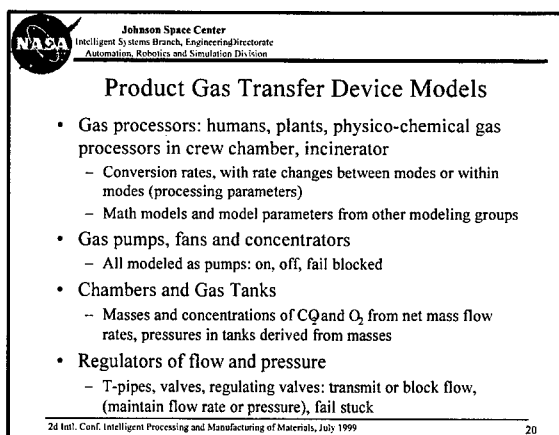
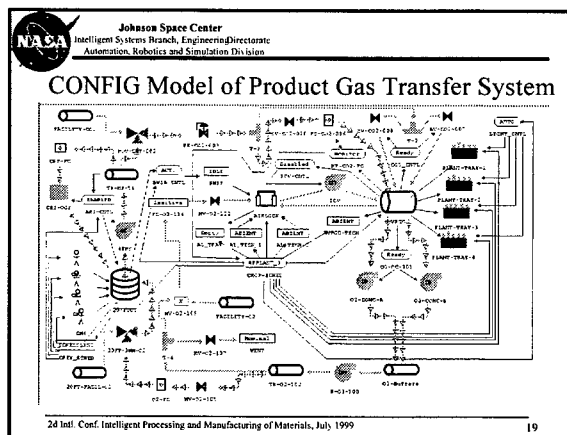
2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 17


 **Johnson Space Center**
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

CONFIG Object-Oriented Models

- The models should be as simple possible and as complicated as necessary for the software testing
- Models of behavior of both components and operations (activities: control, procedures, scenarios)
- Efficient modeling and simulation of components in reconfigurable paths of flow and effort
 - Models recomposed during simulation as direction and activation of interconnections changes
- Reuse and refinement of models and model parts
 - Interactive graphical model building
 - Modular object-oriented libraries
 - Reuse of quantitative models developed for performance analysis

2d Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 18




 **Johnson Space Center**
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

Test Results

- Product gas transfer RAP sequencer performance was evaluated against each software requirement
 - Few errors, quickly fixed
 - More requirements issues than bugs/errors
 - Issues and errors appeared in system-level interactions
- Model deficiencies were mostly due to changes in requirements and design
 - Model was changed incrementally and quickly fixed


24 Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 25

 **Johnson Space Center**
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

System-Level Issues Raised during Testing

- Adjusting transfer rates and setpoints
 - Sequence of control regimes needed to manage CQ after incineration required manual intervention at local controllers
- Reconfiguring transfer hardware and gas supplies (redirecting flow)
 - IMC accumulator venting insufficient to manage excess CO_2 in interacting tanks, would require manual intervention
 - Problem in handover to backup software: failure to switch to facility CO_2 supply to resolve a CQ accumulator problem
 - Gradually dropping pressure in accumulator compensated at first by increased flow in flow-controlled transfer to plant chamber
 - Drop in CO_2 density in plant chamber to alarm level led to switch to backup injection system, but without switching to alternate source

24 Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 26

 **Johnson Space Center**
Intelligent Systems Branch, Engineering Directorate
Automation, Robotics and Simulation Division

Conclusions and Future Work

- The hybrid modeling approach and modeling environment supported the demands of the validation test
 - The full range of model representations was necessary for the test, especially for control and reconfiguration
- Simulation-based testing found problems with software and requirements that would be difficult to detect in more conventional testing, especially system-wide interactions
 - New CONFIG extensions support interactive operator-in-the-loop evaluations, and testing of all three tiers of 3T architecture
- Multiple model types will be needed for engineering and operating autonomous production plants on Mars
 - More detailed models are being developed for physico-chemical gas processing subsystems, for life support and propellant production

24 Intl. Conf. Intelligent Processing and Manufacturing of Materials, July 1999 27

Discrete Modeling via Function Approximation Methods - Towards Bridging Atomic- and Micro-Scales

A. G. Jackson and M. D. Benedict

Materials and Manufacturing Directorate,
Air Force Research Laboratory,
Wright-Patterson AFB, OH

Discrete modeling of processes at the atomic-scale affords practical approaches to complex materials of interest commercially and to the United States Air Force. Reductions in computation times can be large, suggesting the possibility of real-time modeling of thin film growth and the consequent development of processing routes to achieve specific physical and chemical properties. Formulation of the model to be used is critical in achieving such computational gains. Frameworks for these models such as Monte Carlo and Molecular Dynamics can be used conceptually, but they cannot be applied in practice because of the high number of required computations per time step.

The simplest discrete model involves the Potts Model to simulate energies, then to create a partition function of probabilities for various states and configurations, followed by a decision algorithm that determines the state of surface atoms. Although the inclusion of defects, dopants, atom complexes, surface reconstruction and crystal orientations can be included directly in this modeling approach, the resulting collection of behaviors is very entangled with logical and mathematical functions. Hence, the time to exercise the model increases noticeably.

This problem can be reduced dramatically by employing neuro-computing methods. Because neural nets can be trained to represent very complex non-linear relations, substituting neural nets for those methods enables a thousand-times speed-up in atomic-scale simulation. These gains in speed result in near-real-time and real-time atomic-scale models for large numbers of atoms ($>10^4$) on desktop computers enable effective process design and development for thin films and other coatings produced by vapor and liquid deposition techniques. In addition, it is noted that the scale up to micro-scale is conceptually approachable via modification of the discrete space used in the model, opening up a direct connection with existing continuum models.

**DISCRETE MODELING VIA
FUNCTION APPROXIMATION METHODS -
TOWARDS
BRIDGING ATOMIC- AND MICRO-SCALES**

A. G. Jackson* and M. D. Benedict*

*Materials & Manufacturing Directorate,
Air Force Research Laboratory
Wright-Patterson AFB, OH, USA*

IPMM-99

*Second International Conference on
Intelligent Processing and Manufacturing of Materials
July 10-15, 1999, Honolulu, Hawaii
Hybrid Modeling Symposium*

*AvXm Partnership, Dayton, OH

Research Objectives

- **Explore methods for nanoscale materials-process design that are extendable to larger scales**

Research Focus

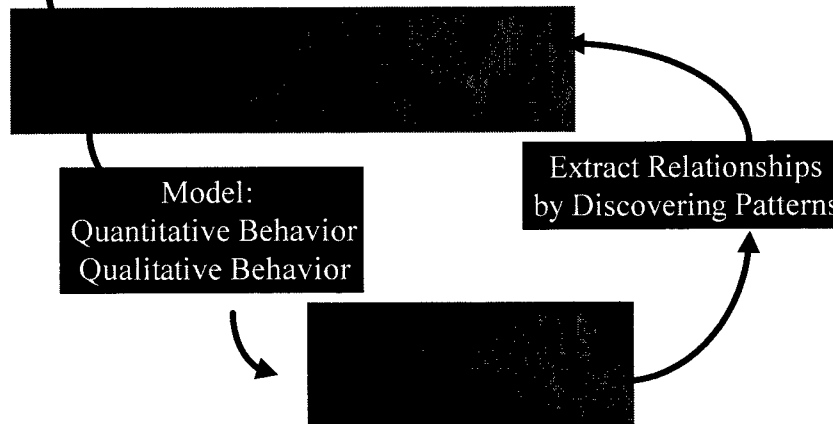
- *Develop* more computationally tractable methods for simulating thin-film materials process design and visualizing resultant crystal structures and the causal processing conditions and phenomena therein,
- *Investigate* the use of patterns: rules/constraints, to simulate thin-film physio-chemical formation and growth via finite state machines,
- *Develop* cellular automata-like methods to improve the speed of simulating micron to millimeter thick films over topographically irregular surfaces for conformal functionality

Atomic Scale Design

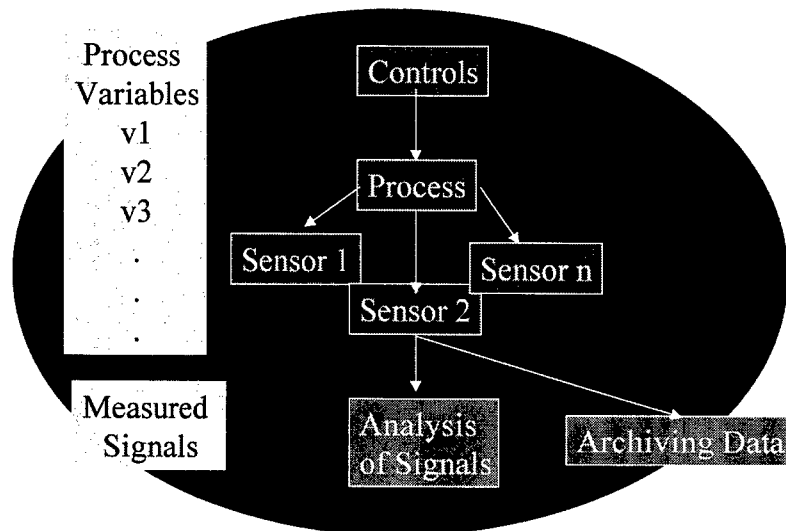
Materials System: Thin Film
Optical, Thermal, Electrical,
Magnetic, Mechanical Properties

Model:
Quantitative Behavior
Qualitative Behavior

Extract Relationships
by Discovering Patterns



Information Flow for a Process



- **Chemical Vapor Deposition**

LaAlO₃ interface coating on Al₂O₃ fibers

More Affordable Ceramic Composites
for HighTemp Engine Components,
e.g., F-16 engine nozzle flaps upgrade

- **Pulsed Laser Deposition**

Superconducting (YBCO) Thin-films

Process Adaptability for Optimized
Thick, Thin & Uniform High T_c films,
e.g., Improved missile phased array radar
in terms of range and resolution

- **Pulsed Laser Deposition**

HARD (DLC) Coatings

Process Adaptability for Hard - Tough
Multi-layer coatings,
e.g., Improved life of Global Positioning Satellites

- **Molecular Beam Epitaxy**

Semiconducting (III-V and II-VI) Thin-films

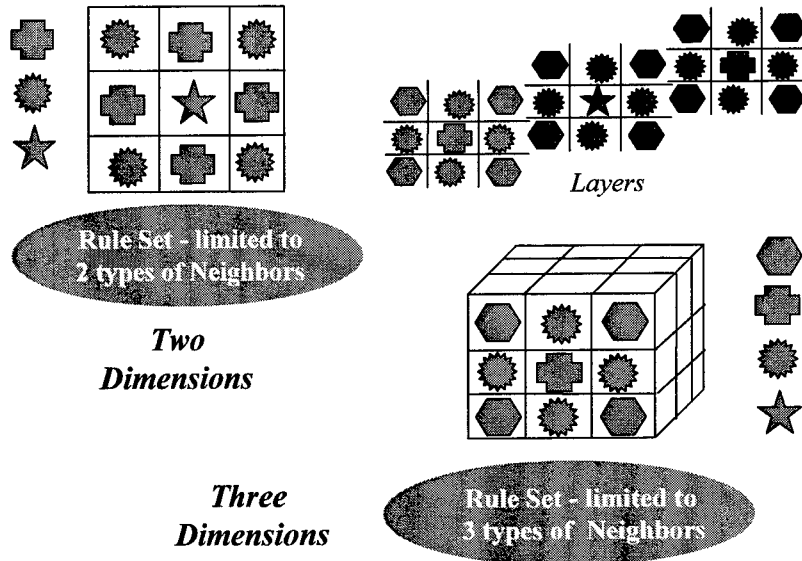
Repeatability for Superlattices & Rugates,
e.g., 2D / 3D opto-electronics
(C³ modulars, sensors, etc.)

Atomic Scale Thin Film Growth

Thermal Expansion
Thermal Vibrations
Surface Diffusion

Structure
Physical Properties
Growth Mechanisms
Processing Conditions

Cellular Automaton



CA toward CA-like Representations

Geometric arrangement of cells
Uniform rule sets
across all neighborhoods

Variations:

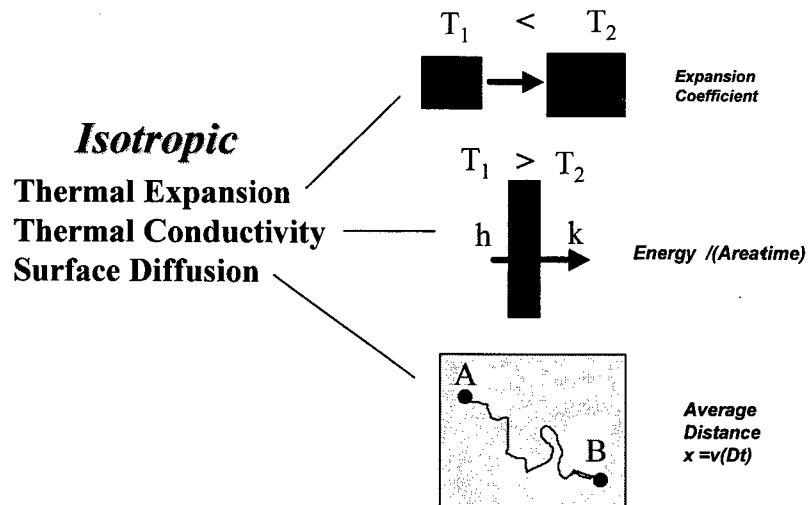
Variable cell geometries

Evolving rule sets

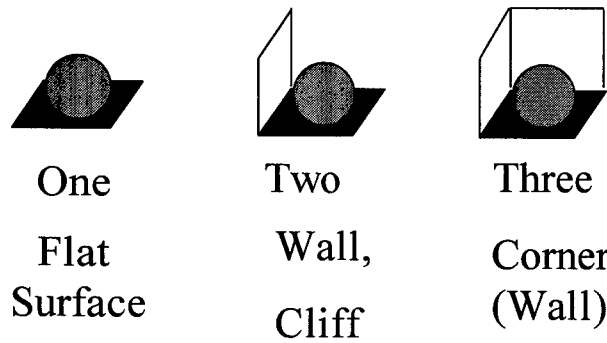
Neighborhood rule sets

CA becomes State Change Algorithm

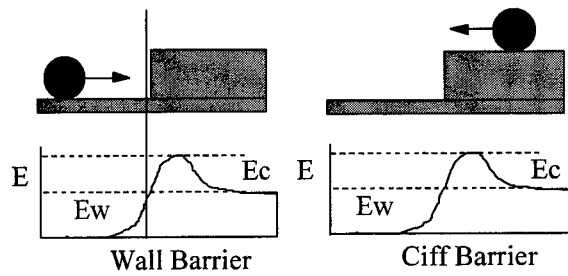
Physical Model



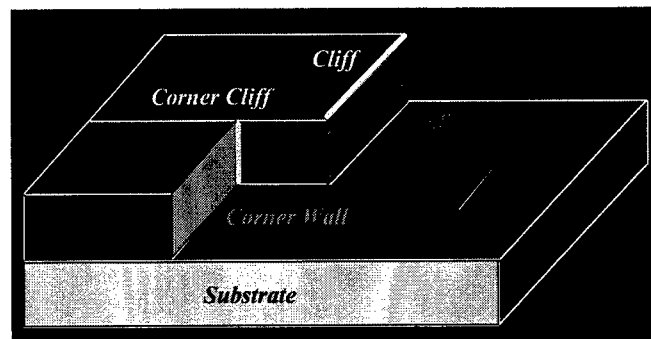
Adsorption Bonds



Simplified Interaction Model



Substrate Surface Model



Computational vs CA-NN

Computational Methods

Constants
f1, f2, f3, ...
Integrations
Differentiations
Equations of Motion
 ...

*Problem is
 Computationally
 Demanding*

CA-NN Method

Generate Exemplars
Train the Net Off-line
Predict On-line

Problem is decomposed
Exemplar Generation
Training
Fast Prediction

Decision Algorithm to Choose State

Conditions

*Probabilities
 of Particular
 States*
*Physical
 Parameters*
*Neighborhood
 Configuration
 and State*

Engine

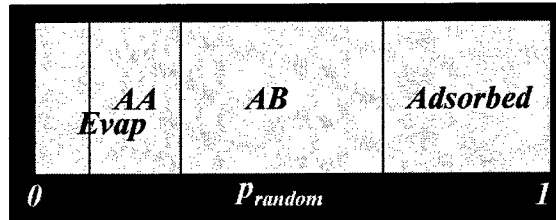
CA-NN
Training Set
Testing Set

Decision

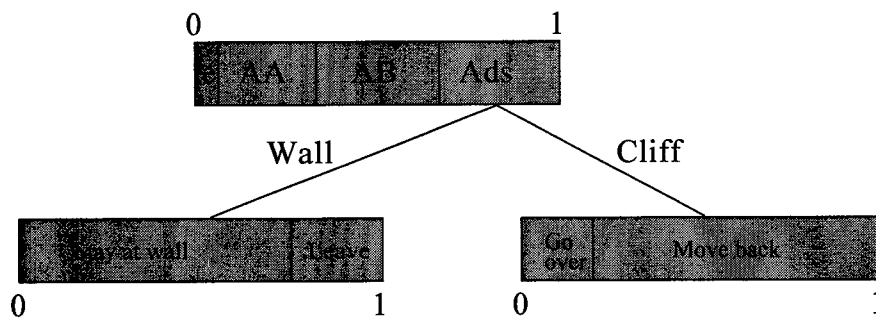
*Output
 State*



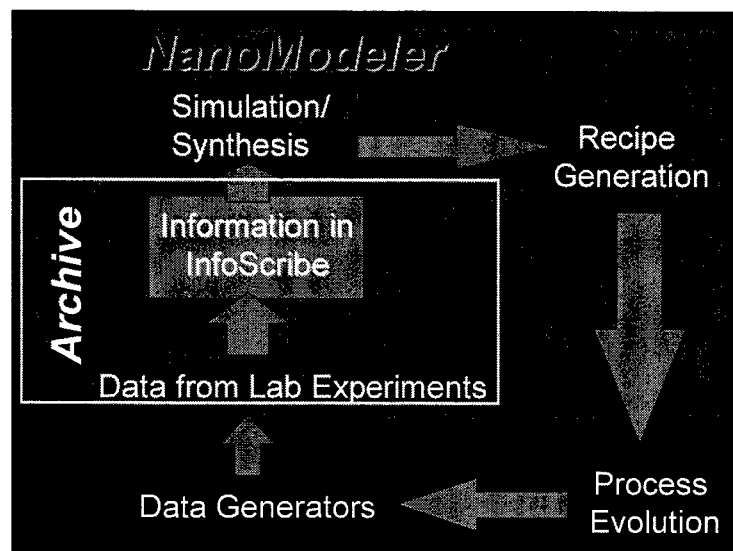
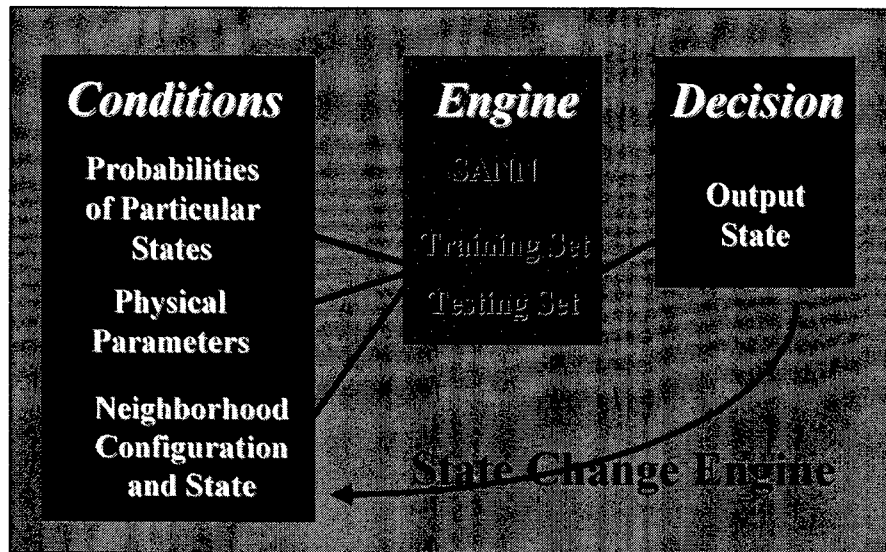
Barrier Probabilities Diagram



Barrier Probabilities Diagram with Wall and Cliff Added



Implementing the Model



Future Research

- Qualitative Models
- Pattern Mining
- Design Environment
- Computation Simplification
- Physical Models at Various Time Scales



Surface Microstructure Predictions from Atomistic Rule Set Cellular Automata

M.O. Zacate, K.J.W. Atkinson, R.W. Grimes and P.D. Lee

Imperial College, London University, London, England, UK

When a specific microstructure is required, if the preparation variables increase beyond a few, it is very difficult and expensive to determine the optimum conditions experimentally. Consequently there is considerable interest in predicting conditions via computer simulations. Since ultimately, microstructure depends on processes occurring at the atomistic level, to be fully transferable, it is desirable that such a model is atomistically-based. This should also allow us to include the role of all types of chemical and crystallographic defects explicitly.

In this study, we begin by calculating the energetics associated with the way in which individual gas atoms interact with a specific metal surface. Both perfect and defective metal surfaces are considered. The energetics are translated into rule sets which form the basis of the cellular automata. The rule sets involve both temperature and gas atom flux as variables. The result is a model which can quickly, explicitly describe the evolution of 104 surface sites over 10^{-6} seconds with very modest computing facilities.

In the simulations, the formation and growth of domains which exhibit critical behavior are observed. That is, the rate of growth is not a well-behaved function of temperature or flux but exhibits a region in which the rate of growth suddenly falls to zero. Surface defects are also predicted and have a dramatic effect on growth rates.

Presentation Outline

Part 1. Atomic scale simulation methodology.

Figure 1.1 The classic length scale view of modelling techniques.

This is used to show exactly where this presentation sits in terms of the modelling hierarchy and what we intend to achieve in general modelling terms.

Figure 1.2 How modelling at the atomic level can be approached.

Figure 1.3 How to calculate forces/energies.

Figure 1.4 Relationship of computer simulation to experiment and basic theory.

Part 2. Surfaces and their structures.

This part of the presentation outlines the problems faced when attempting to model the surfaces real materials.

Figure 2.1 Atomic structure of a (111) surface of ZnCr_2O_4 .

Shows the significance of surface relaxation, that is, real surfaces do not look like a piece of perfect cleaved crystal. Simulation carried out using energy minimisation.

Figure 2.2 Classification of surfaces.

Note the fact that a type III surface is not stable because it possesses a dipole moment perpendicular to the surface.

Figure 2.3 Stabilisation of a type III surface by the introduction of surface point defects.

Figure 2.4 Configurations of point defects on the (100) surface of UO_2 .

There are two ways of arranging the two oxygen ions over four sites. When a 2×2 unit cell is considered the number of possibilities increases to 105.

Figures 2.5 & 2.6 Molecular dynamics simulation of a 3656 atom cluster of CaF_2 at 300K viewed from two angles.

These show the cluster initially before the MD run and after 110 ps which corresponds to over 10^5 time steps, clearly a significant calculation. Thus modelling of large surfaces using this methodology is not feasible - however, it is important and useful for verification and for test cases.

Figure 2.7 HREM image of commercial CeO₂ powder, the arrows show the presence of surface facets.

The previous overheads described the atomic structure of specific surfaces (i.e. the atomic structure of a flat surface). This figure shows that in addition to atomic relaxation and point defects, even on the nano scale, i.e. on a level just beyond the atomic level, real surfaces are not perfectly flat but have ledges and facets. This should also be taken into account when modelling a real material.

Figure 2.8 What we need to be able to be able to model on the atomic level.

Finally, to model surface evolution at the microstructural level we must address possible the computationally most demanding problem. We need to model the behaviour of surfaces containing many thousands or even hundreds of thousands of atoms.

Part 3. Microstructure predictions from atomistic rule set cellular automata - a solution to the length scale problem for complex surfaces.

Most of the results and methodology presented here are available as electronic reports on the following web site:

<http://abulafia.mt.ic.ac.uk/USAF>

<http://abulafia.mt.ic.ac.uk/USAF2>,

or have been published in the open literature:

Zacate M. O., Grimes R. W., P. D. Lee, S. R. LeClair & A. G. Jackson "Cellular Automata Model for the Evolution of Inert Gas Monolayers on a Calcium (111) Surface" *Modelling Simul. Mater. Sci. Eng.*, in press.

Figure 3.1 Methodology summary.

Figure 3.2 Definition of the crystallography of the problem.

Gas atoms may either occupy B or C interstitial sites on the surface. If a B site is occupied, the adjacent C sites cannot also be occupied due to steric effects.

Figure 3.3 Formation of a growth fault.

This shows a 2D equivalent to a stacking fault. The area around the fault has a lower density of has atoms. This is due to the crystallography through the steric effect.

Figure 3.4 Results for the microstructural evolution of a layer of Ar on a 200x200 atom surface of Ca after 2,000 time steps (i.e. $\sim 10^{-6}$ s). Six temperatures are shown. Note the strong, non-linear temperature dependence.

Figure 3.5 Attachment energies of Ar atoms to the surface adjacent to a growth fault.

Note the lower energy of atoms close to the fault and the rather different energies depending on the details of the local geometry i.e. the number of nearest neighbour gas atoms which depends on the geometry of the growth fault.

Figure 3.6 Atomic scale explanation for why during the CA simulations we observe that small curvature diminish. That is, these atomic scale simulations reproduce the macroscopic effect of line tension.

Figure 3.7 Argon convergent coverage v's temperature for six different incident argon fluxes.

Note the critical behaviour, particularly for small incident flux values.

Figure 3.8 Average domain size (in % of surface area) v's temperature for six incident argon flux values after 95 time steps.

Figure 3.9 Average domain size (in % of surface area) v's temperature for six incident argon flux values after 995 time steps.

Note how the average domain size increases as a function of temperature along a curve that is independent of flux as long as the temperature is less than some critical domain temperature.

Figure 3.10 Domain growth on a rough surface as a function of time.

This 100x100 site surface has 4% coverage of additional Ca ions randomly distributed in the over layer. This rough surface is then subjected to an 8% argon flux at 90 K. Results are shown after various times. Note how in the last two, the domain size has not changed despite a time difference of an order of magnitude. the microstructure has been pinned by the roughness.

Figure 3.11 Domain growth on a rough surface as a function of roughness.

The microstructure is now shown after 30,000 time steps but with a variety of different roughness factors to show how domain size depends on the degree of roughness. Of course we must be aware of what is meant in this case by roughness - a random distribution. By July we will have data concerning other types of roughness.

Surface Microstructure Predictions From Atomistic Rule Set Cellular Automata

M. O. Zacate, K. A. J. Atkinson, R. W. Grimes and P. D. Lee

Department of Materials, Imperial College, London, UK

Introductory remarks – motivation.

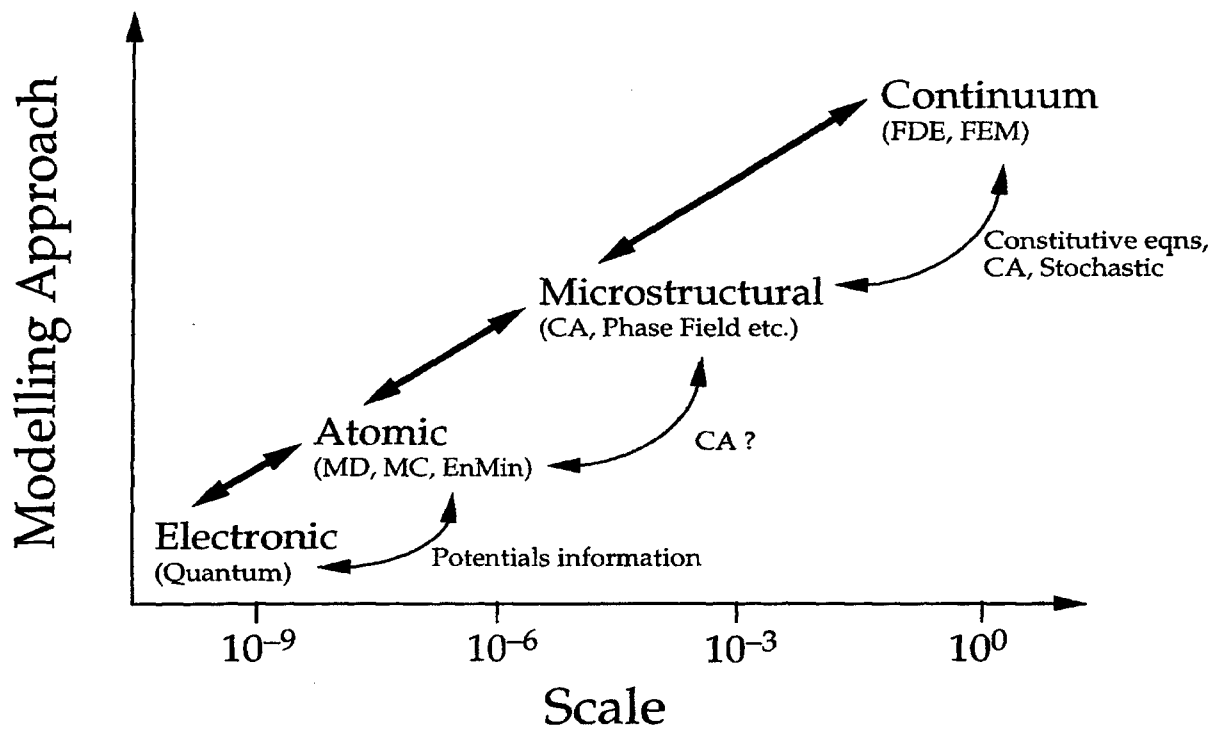
It is now computationally possible to bridge the simulation length scale gap from atomistic to microstructure. The aim of such an endeavour is clear: to be able to predict how processes which occur on the atomic level affect and control microstructural evolution. In addition, it will be critical to translate changes to variables which effect the atomic processes into changes in process control parameters. Some of these, such as temperature are obvious, others will be more challenging. As a first step, it is our intention to discover or better understand the competing mechanisms which govern the connections between atomic processes and microstructural evolution and explore to what extent they are universal or system specific.

To appreciate the issues involved we will begin by considering a familiar example. In CVD or PLD, surfaces may grow via a 2D layer-by-layer process, forming a smooth morphology, or by a 3D column process forming a rough surface. 2D growth results in crystallites bound by well defined specific planes. 3D growth may result in almost spherical grains independent of the surface crystallography. The underlying mechanisms that are responsible for the change from 2D to 3D growth regimes are influenced by substrate and vapour stream temperatures, vapour stream and environmental pressures, deposition rate and profile (e.g. pulsed or continuous) and the geometry of the deposition chamber and deposition process. These are machine variables and form a unique n dimensional set for each different combination of film chemistry and substrate. Nevertheless, certain broad characteristics of growth are similar from system to system. For example, a smooth to rough growth transition may be caused by changes to the machine variables that affect the rates of atomic processes such as surface diffusion, absorption and desorption rates, although the characteristic energies of these processes depend on the chemistry of the specific system and local crystallographic geometry. The complexity of an n dimensional set of machine variables can therefore potentially be translated into the complexity of atomic processes. The resulting variety of atomic

processes must all be modelled because the most common atomic process will probably not be the deciding factor i.e. the rate determining step. Furthermore, it is not clear that the rate determining step will remain the same over all pressures and temperatures.

In this talk we will present results of very recent studies which have tried to address the length scale problem for a specific model system, rare gas atoms on a metal surface. In this case we are easily able to define all the possible crystallographic environments in terms of a local atomic model, a smooth or locally rough Ca (111) surface. We have then investigated the effect of three model "machine variables", temperature, incident gas stream flux and the roughness of the substrate. The results suggest that the microstructure of the resulting surface will develop in ways that are non-linear to the extent that they may be described as critical phenomena. It is possible that identifying these critical points, as a function of machine variables will become a central theme of process modelling.

Multi-scale modelling



How To Do It !

Static Calculations

Ions which are moved until they reach zero force; minimisation of enthalpy. Quasi-harmonic approximation; average positions, on a time scale which is long with respect to an atomic vibration ($\sim 10^{-13}$ s).

Molecular Dynamic Calculations

Ions have discrete velocities and accelerations determined by solving Newton's equations of motion. Positions are updated iteratively. Atomic vibrations and kinetics modelled explicitly; leads to well defined thermodynamic parameters.

Monte Carlo

Changes to the state of the system are considered on a random basis but excepted subject to energy criteria. This yields information concerning atomic distributions at different temperatures without modelling kinetics explicitly.

Free Energy minimisation

Similar to static calculations except that the vibrational entropy is also included by summing over the vibrational density of states.

How to calculate forces/energies

Quantum Mechanically

Hartree-Fock

Density functional

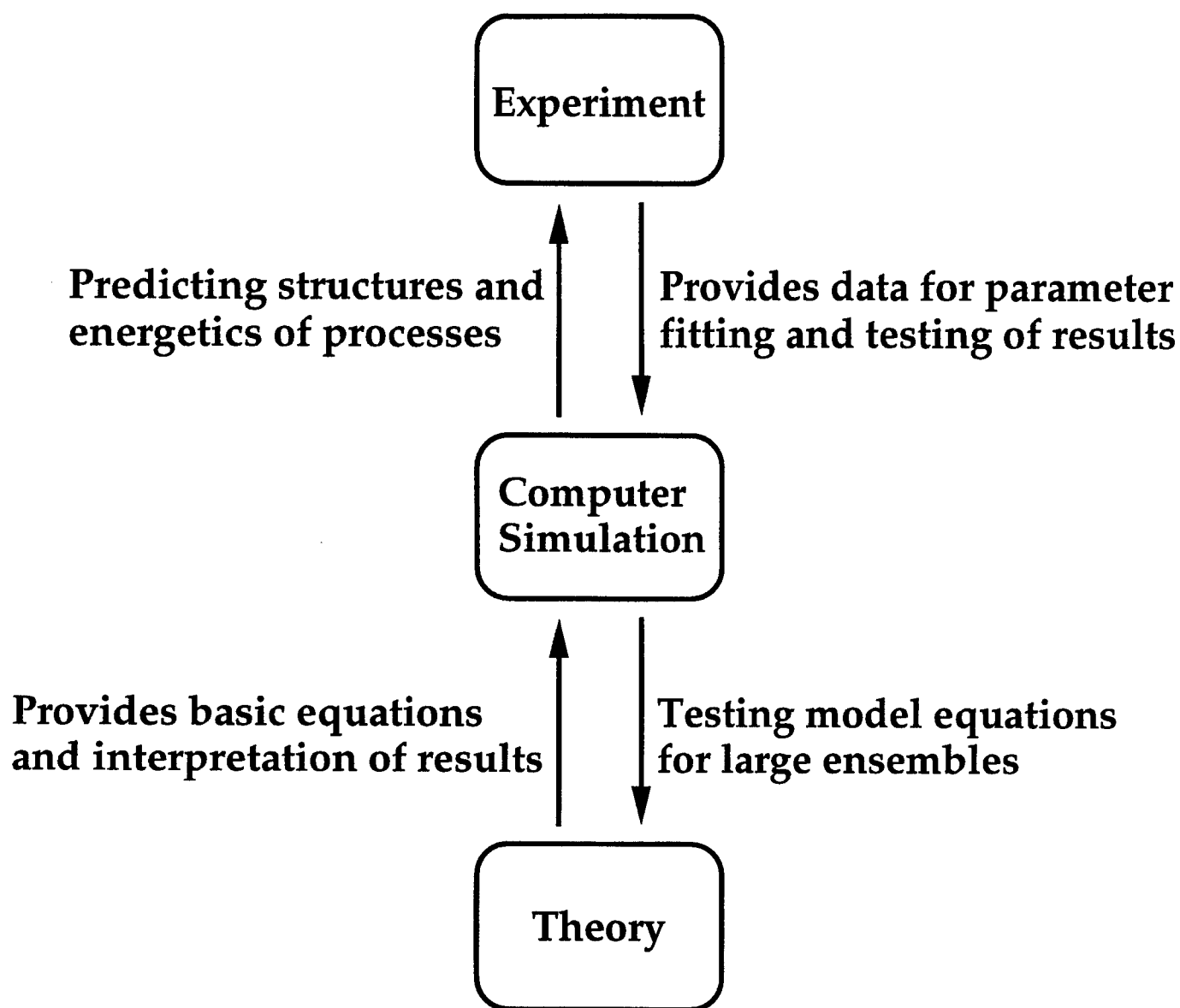
Semi-empirical

Effective potential functions

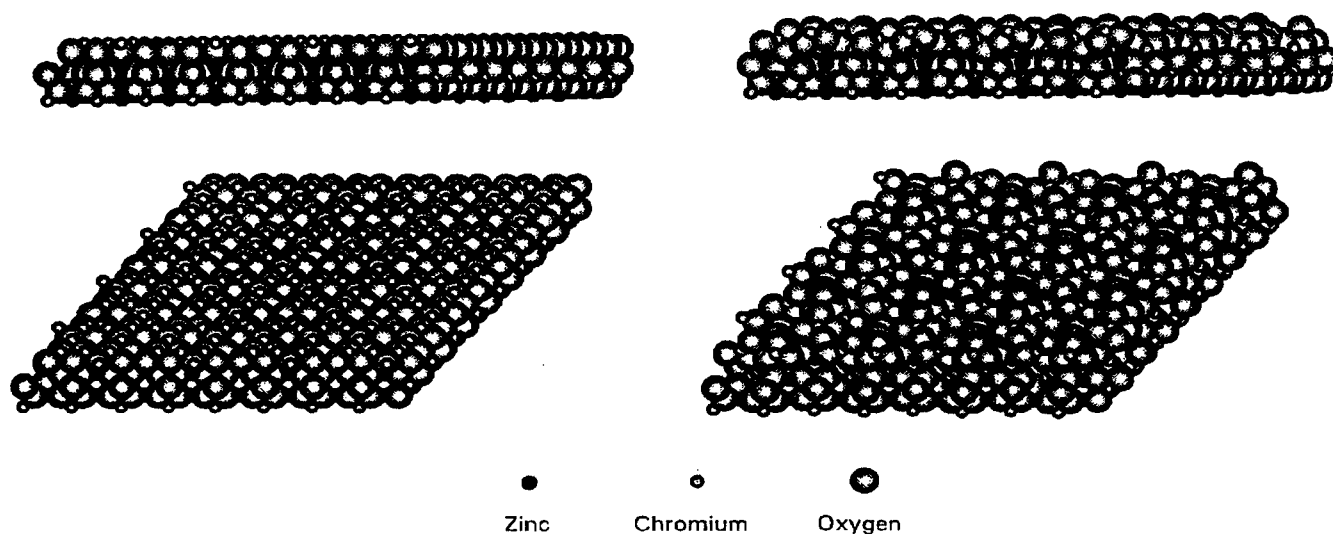
A variety of functional forms.

These are computationally cheap but are usually specific to a system or a series of related systems.

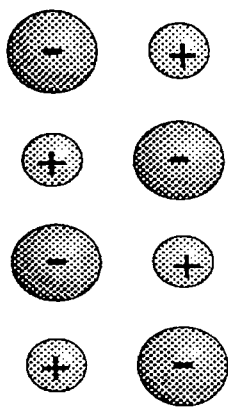
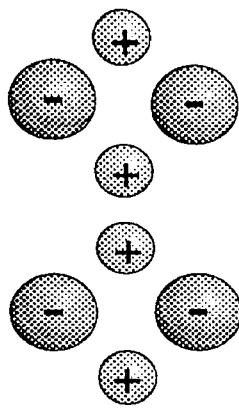
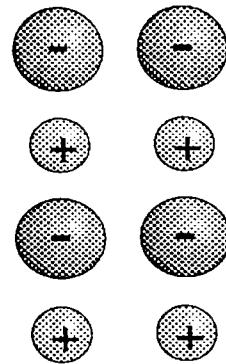
The choice depends on the number of atoms it is necessary to model and the available computational resource.



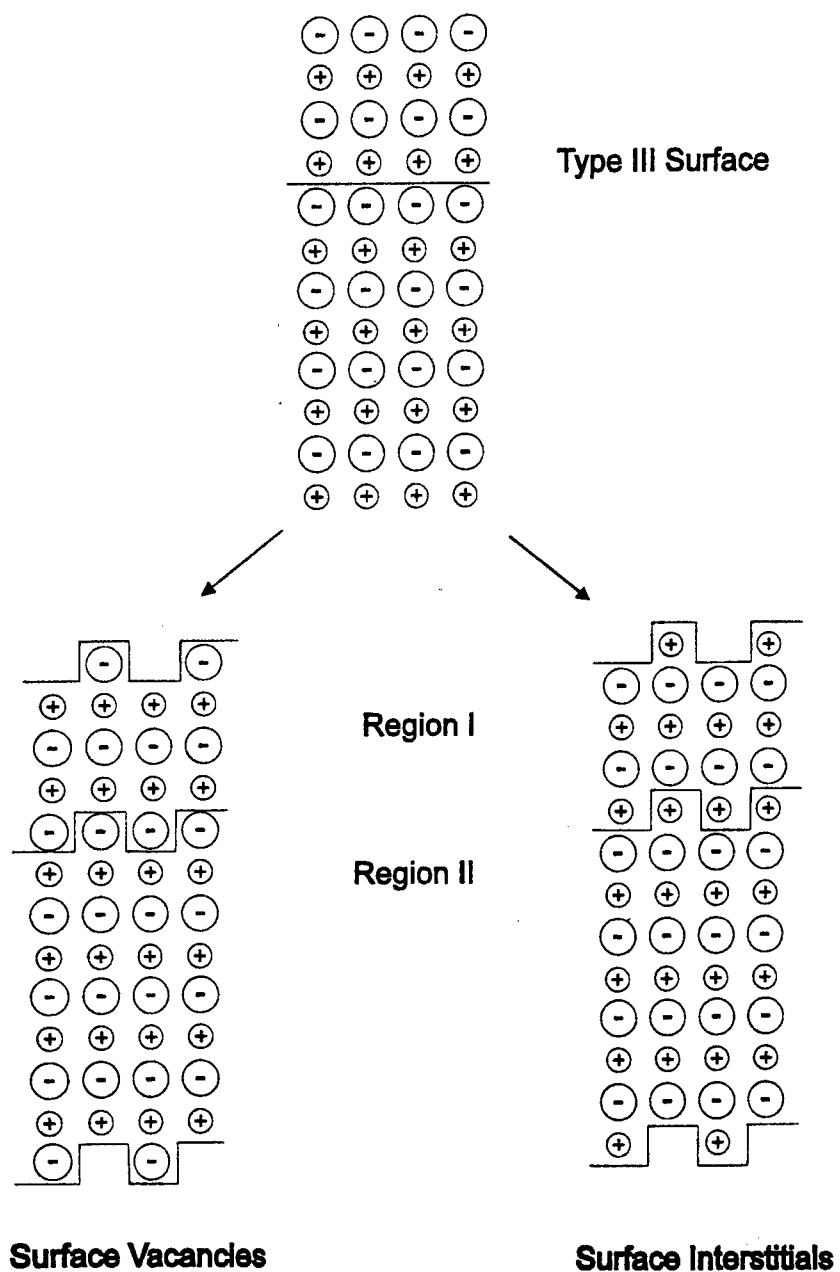
Relationship of computer simulation to experiment and basic theory.



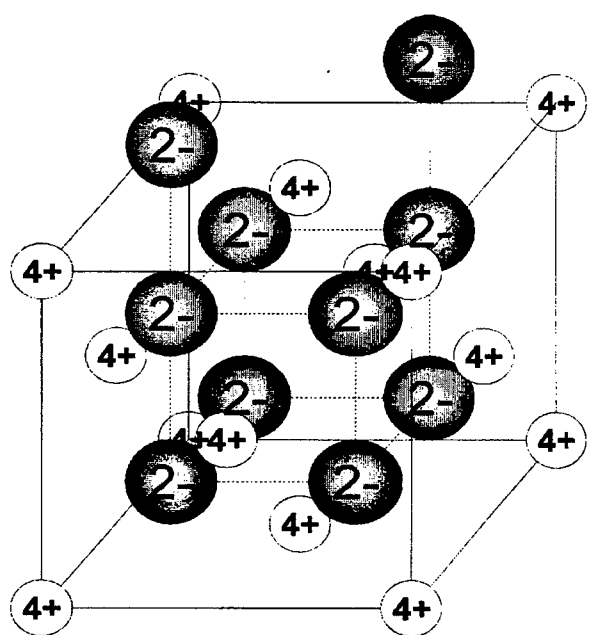
Atomic structure of a (111) surface stabilized by $(6V_{Zn}^{II} + 4V_{Cr}^{III})$.

**Type I****Type II****Type III**

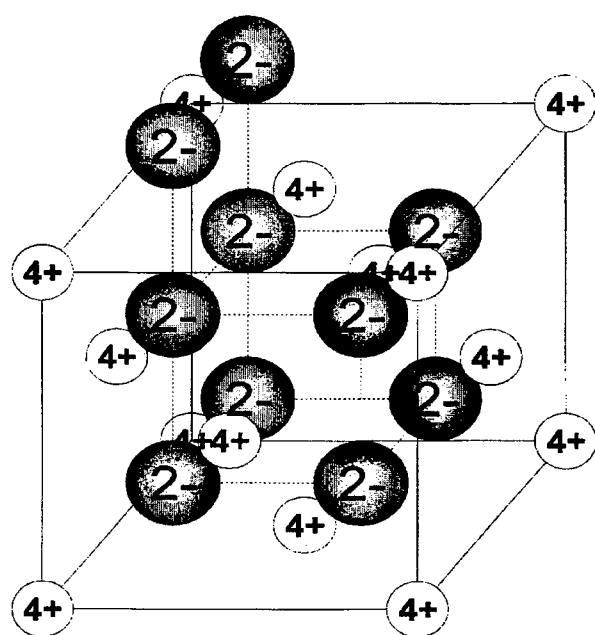
Classification of surfaces



Neutralisation of a dipole on a Type III surface

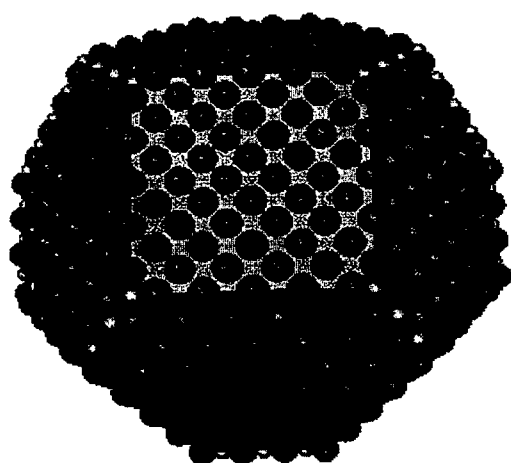


A

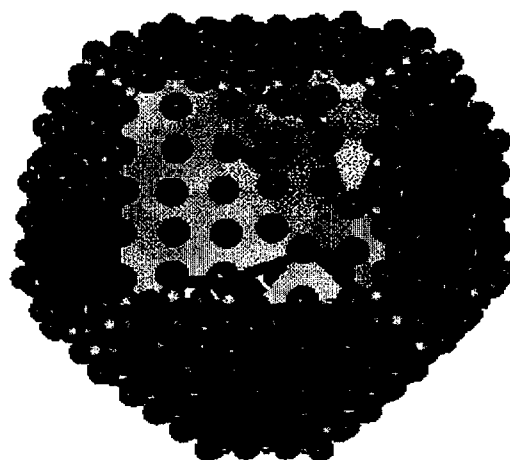


B

A 3656 atom cluster of CaF_2 equilibrated at 300K

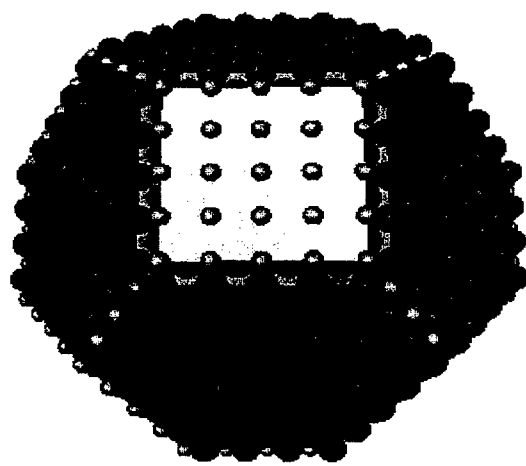


$t = 0$

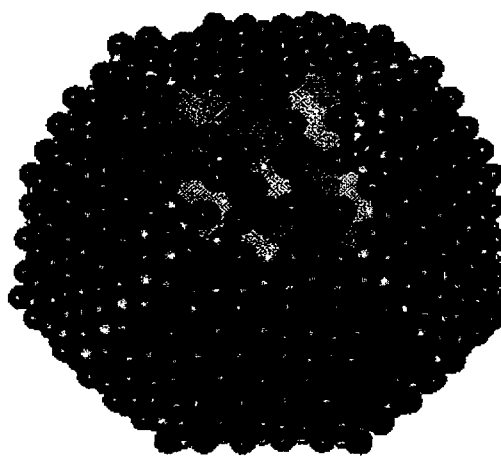


$t = 110 \text{ ps}$

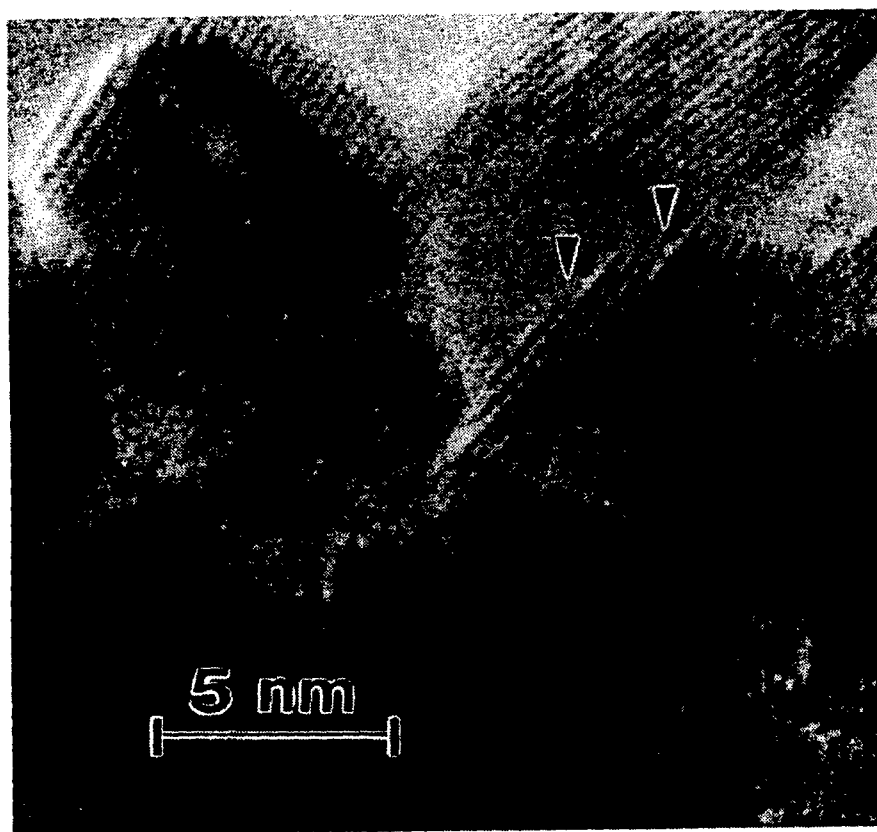
A 3656 atom cluster of CaF_2 equilibrated at 300K



$t = 0$



$t = 110$ ps



Providing energetics for deposition processes.

Cluster formation in the phase above the surface.

Results in an equilibrium distribution

Momentum profile for attachment of molecules to a surface

Depends on the morphology of the surface
Molecules will be different to clusters

Activation energy for molecule migration along a surface

Identification of the rate determining step

Diffusion constant for molecular and cluster surface migration

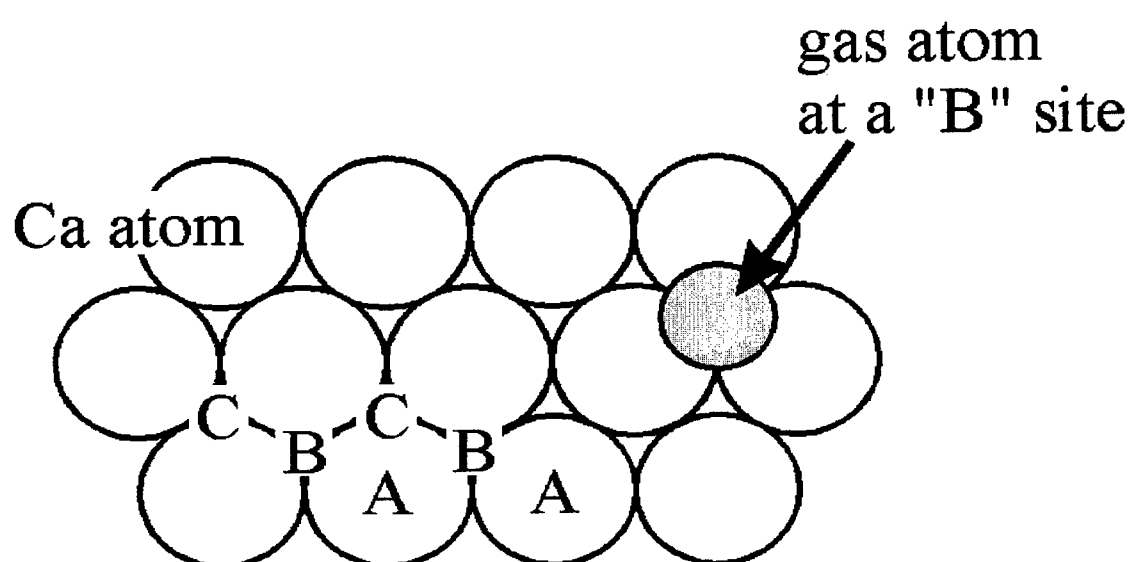
Depends on the morphology of the surface

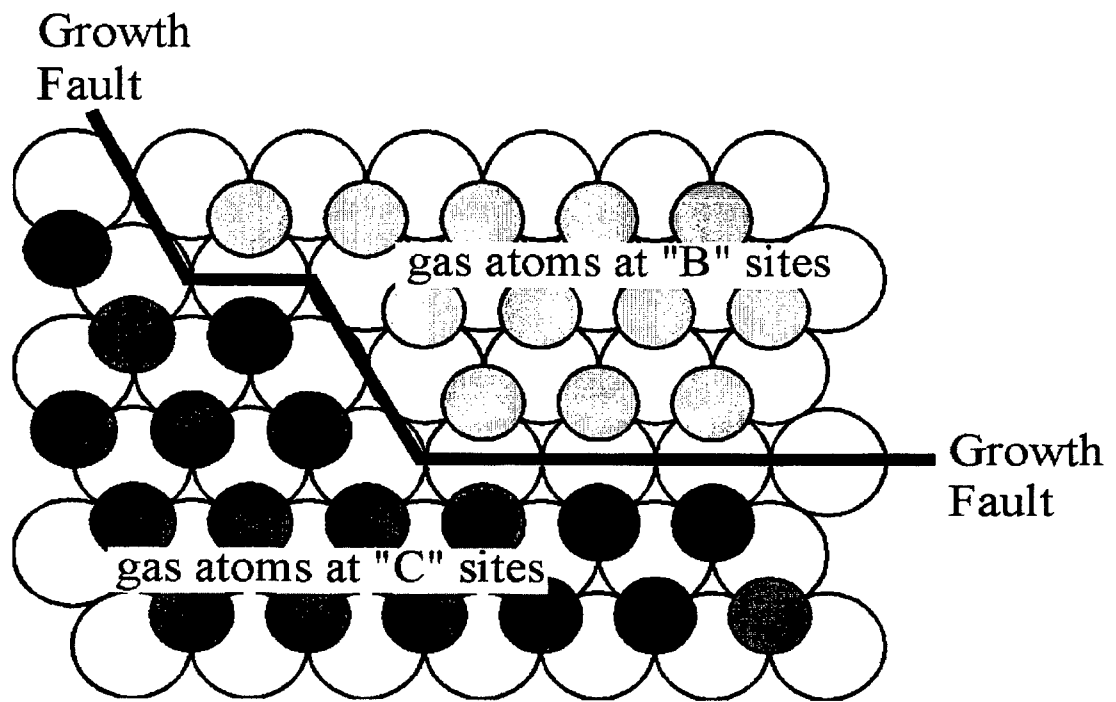
Cluster formation/destruction on a planar surface.

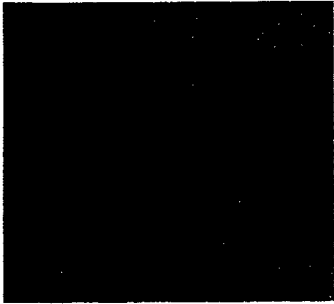
Formation of islands, broken-up by impact events

Attachment of mobile molecules to kinks and ledges

Surface growth mechanisms via kinks



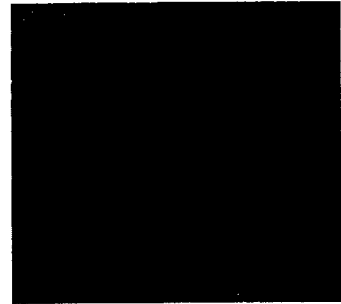




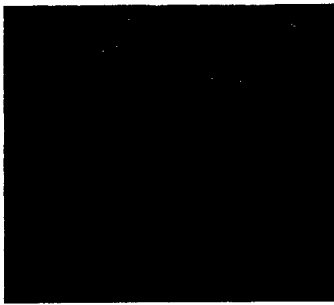
23K



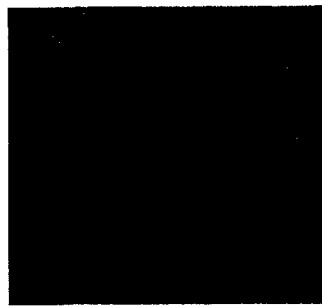
53K



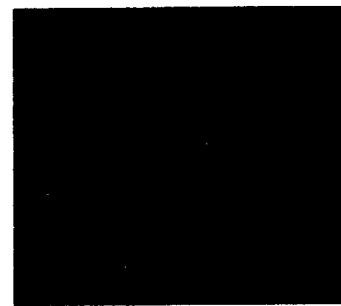
93K



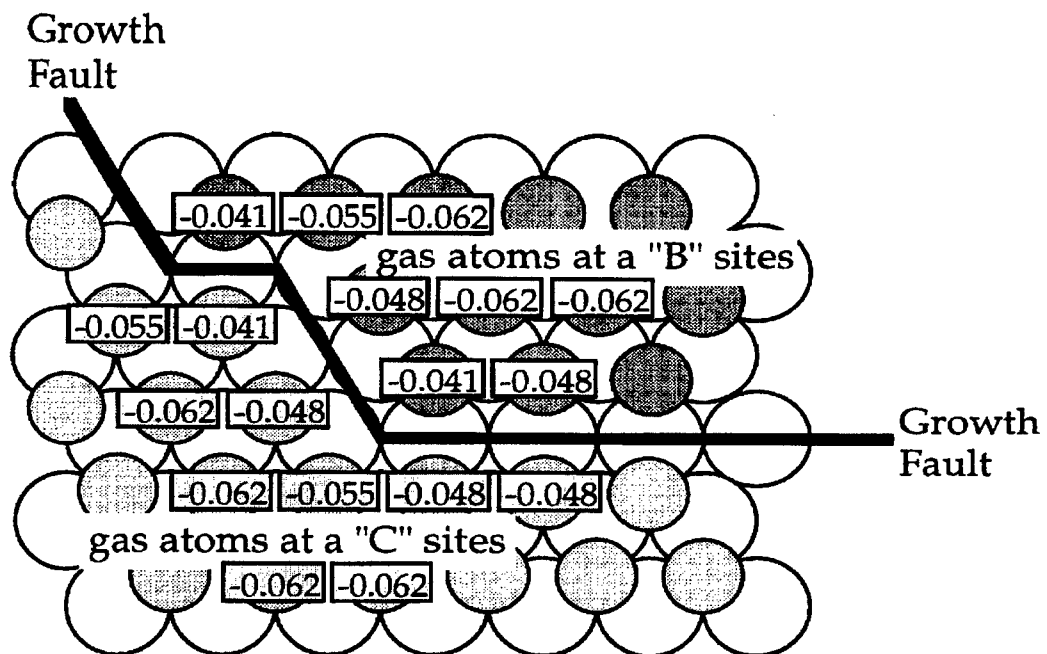
123K

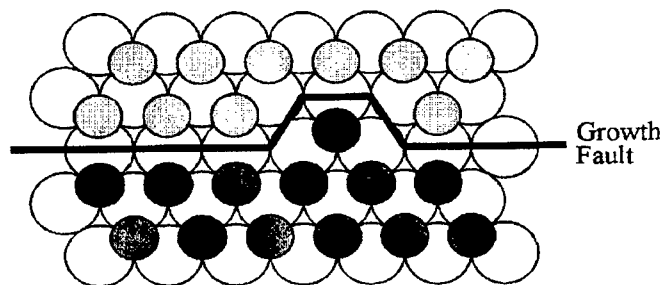
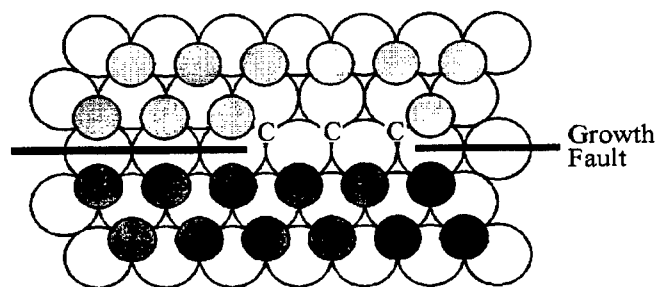
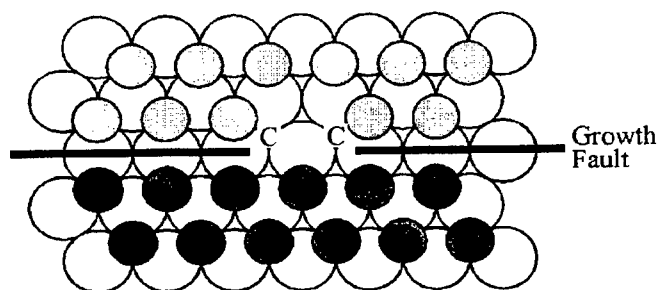
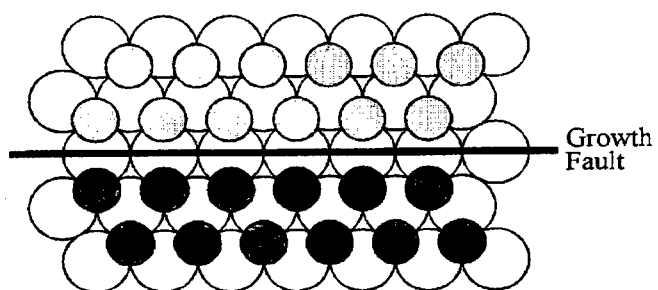


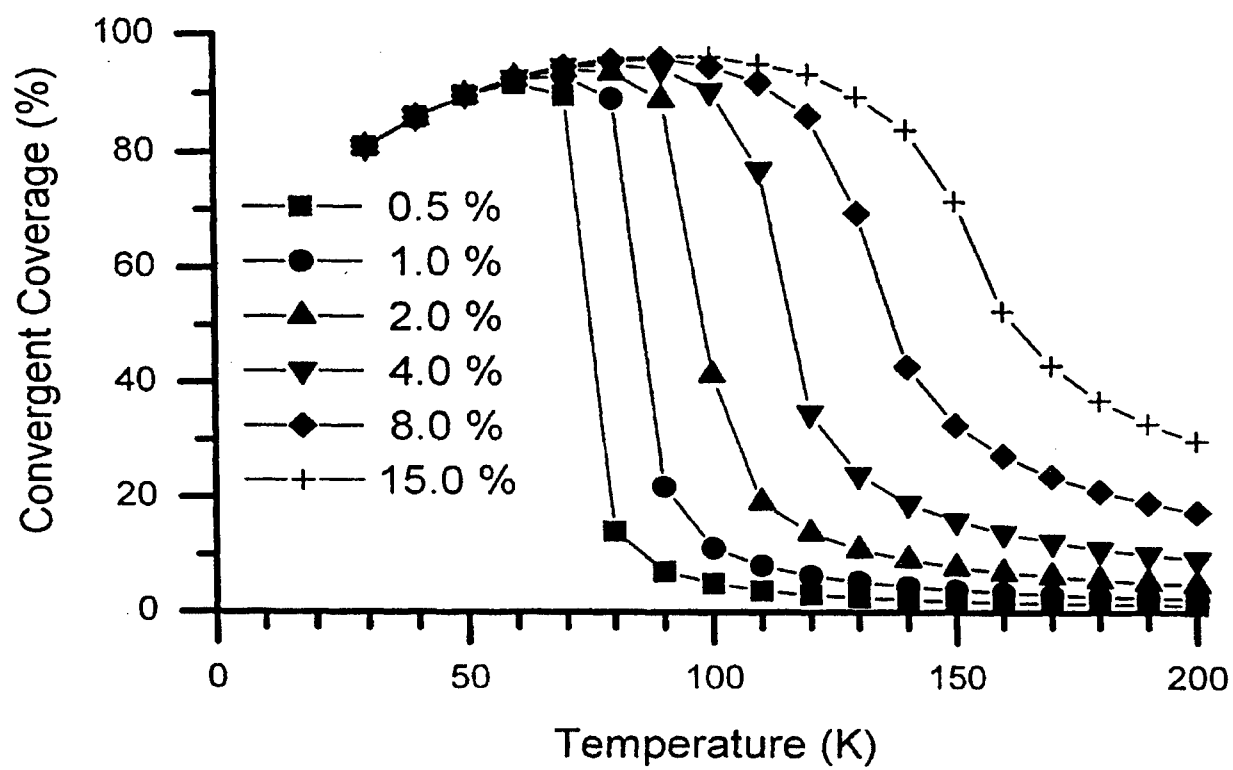
143K

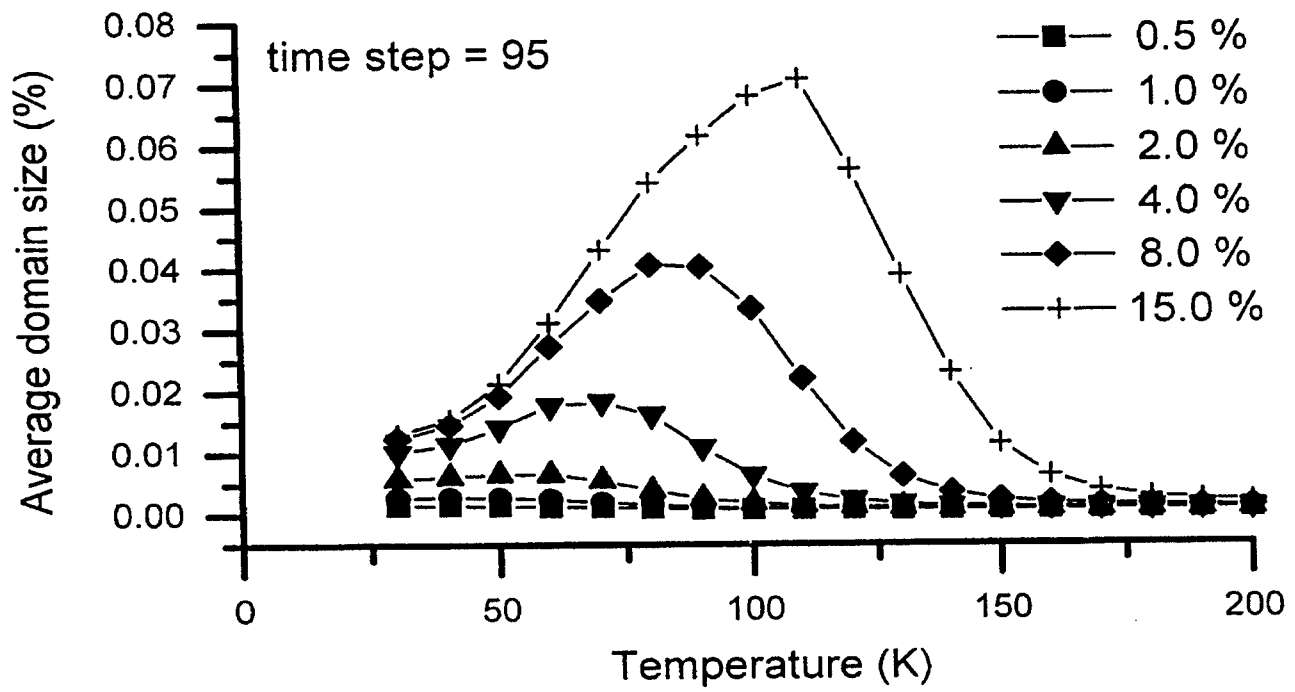


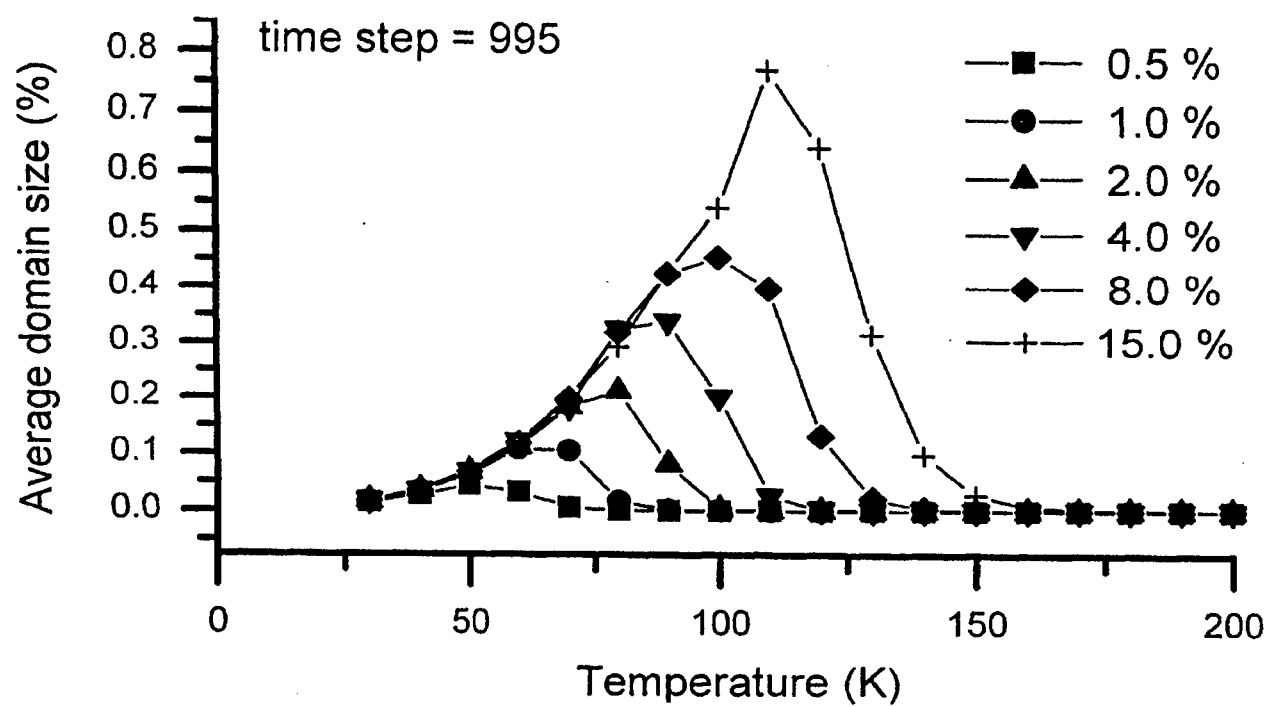
170K

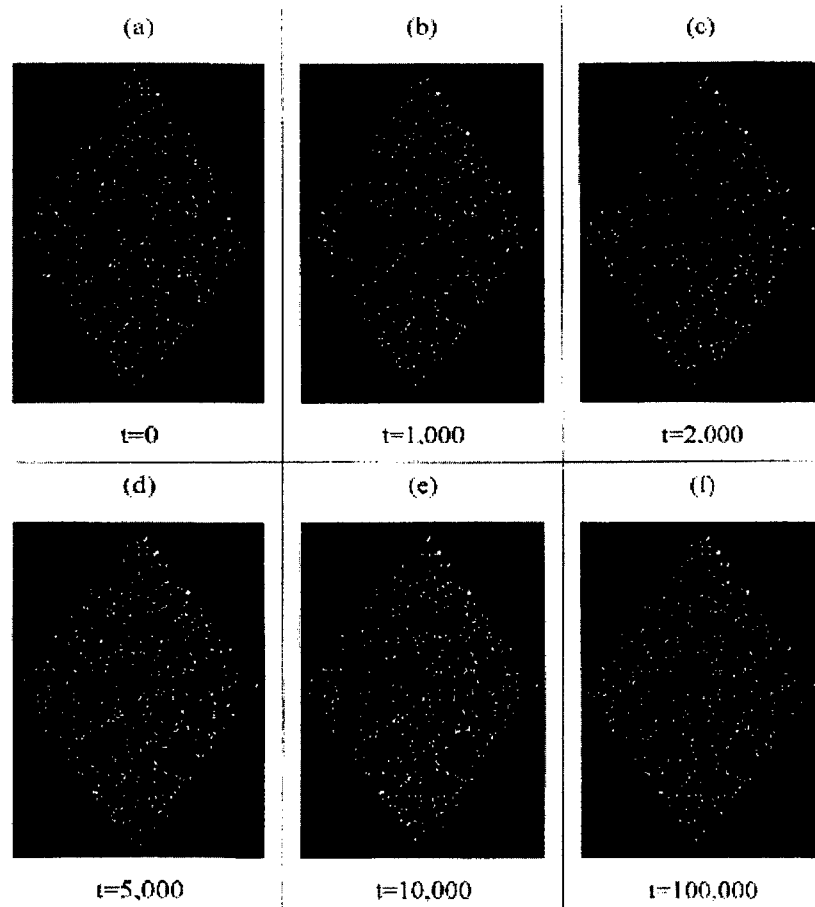


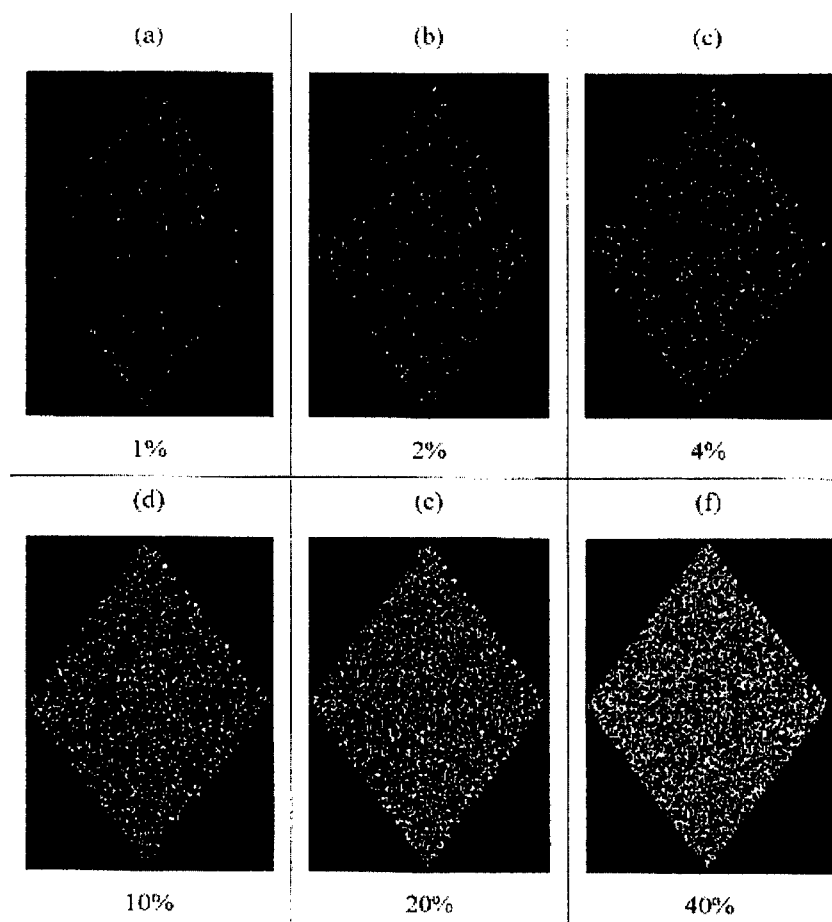












Fuzzy Molecular Modeling

David A. Ress

North Carolina State, North Carolina, USA

The past century has seen a tremendous growth in the knowledge and understanding of chemistry and chemical compounds. This knowledge has enabled researchers to develop new advanced materials such as polymers and composites which have superior characteristics than naturally-occurring materials. Yet, even with the most advanced equipment, we still cannot accurately model chemical properties such as lattice parameters and bond angles in simple tetragonal structures such as the Chalcopyrite family of compounds. The literature is full of researchers reporting lattice parameters and bond angles for chemicals which are, in fact, imprecise. The imprecision results from the analytical techniques and imprecision in the samples and measurement process. The research in this paper presents a new method for representing lattice parameters and bond angles -- through the use of fuzzy logic. Since fuzzy logic by definition is imprecise, it is a natural means to represent the imprecision of lattice parameters and bond angles. The fuzzy lattice parameters are created by collecting parameter values from literature. From the minimum, maximum and average of these parameter values, a fuzzy number is created which represents the imprecision in that parameter. Then, through the use of fuzzy arithmetic operators and recently-developed fuzzy trigonometric functions, atom locations within the unit cell and bond angles can be calculated. Using fuzzy logic in this manner benefits us by allowing the direct representation and calculation of the imprecision found in chemical compounds which arises from measurement error, structural defects, and thermal and vibrational characteristics.

FMM: Fuzzy Molecular Modeling

David Ress

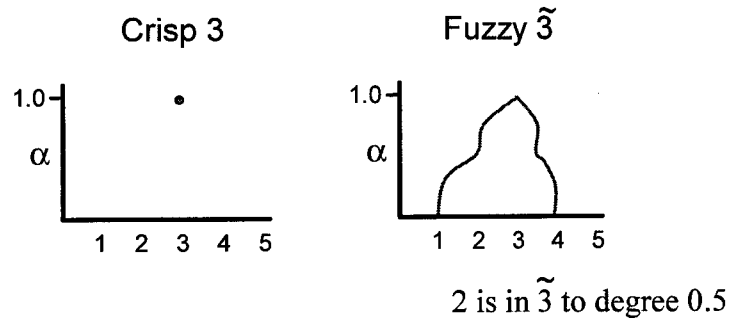
Air Force Research Laboratory
Material Process Design Branch
Wright-Patterson Air Force Base, OH

Why FMM?

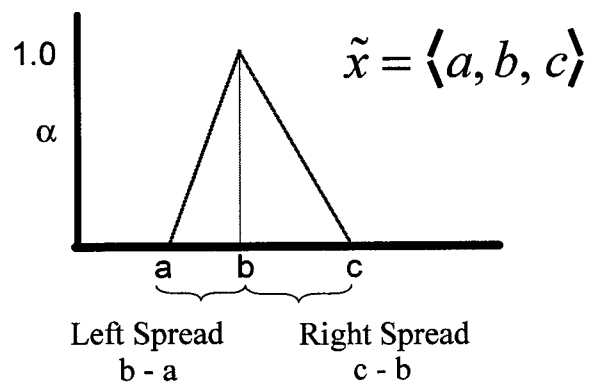
- Difficulties in Molecular Modeling:
 - Measurement Error
 - Structural Defects
 - Thermal Variations
 - Vibrational Effects
- These lead to difficulty in repeatability in measurement, so what is correct?

Fuzzy Set Theory

- Everything is true to a degree.



TFN: Triangular Fuzzy Number



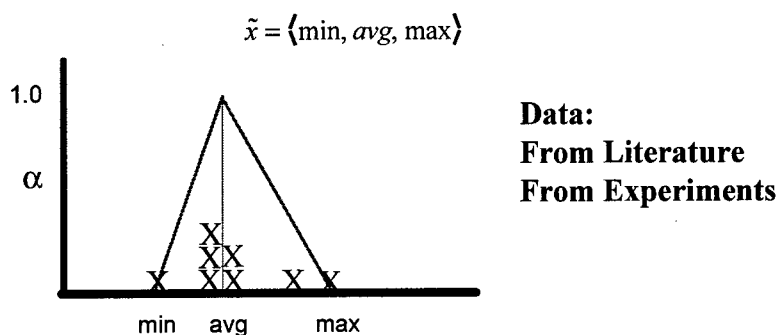
Arithmetics of TFNs

- Just like crisp numbers, TFNs have a rich set of arithmetic operators.
 - Example: $\langle 2, 3, 5 \rangle + \langle 1, 2, 3 \rangle = \langle 3, 5, 8 \rangle$
- Additionally, recent research has developed fuzzy trigonometric functions to support TFNs in design and manufacturing activities.

FMM: Approach

- Collect data to form TFN lattice parameters.
- Create a 'fuzzy' unit cell
- Calculate fuzzy bond lengths and fuzzy bond angles

FMM: Collect Data



This step provides the fuzzy lattice parameters.

FMM: The Fuzzy Unit Cell

- Apply the fuzzy lattice parameters to the Wyckoff coordinates.
- This produces fuzzy atom locations, i.e.,

Fuzzy Lattice Parameters (Å)	Wyckoff Coordinates	Fuzzy Atom Locations (Å)
$\tilde{a} = \langle 5.773, 5.781, 5.785 \rangle$	1/2	$\tilde{a} = \langle 2.887, 2.891, 2.893 \rangle$
$\tilde{b} = \langle 5.773, 5.781, 5.785 \rangle$	1/2	$\tilde{b} = \langle 2.887, 2.891, 2.893 \rangle$
$\tilde{c} = \langle 11.550, 11.602, 11.642 \rangle$	1/2	$\tilde{c} = \langle 5.775, 5.801, 5.821 \rangle$

&



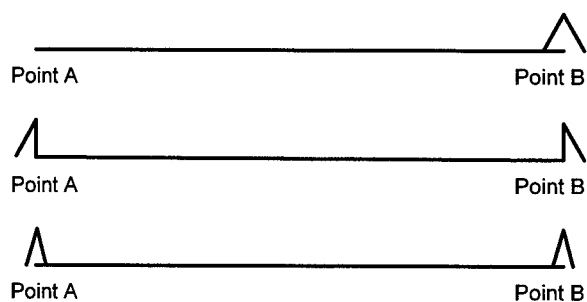
FMM: Fuzzy Bond Angles and Fuzzy Bond Lengths

- With the fuzzy atom locations known, it is possible to construct fuzzy lines (fuzzy bond lengths).
- With fuzzy lines known, calculate the fuzzy bond angles, i.e.,

$$\tilde{\theta} = fa \cos \left(\frac{(\tilde{Dist}_{C1-A6})^2 + (\tilde{Dist}_{C1-B1})^2 - (\tilde{Dist}_{A6-B1})^2}{2 * \tilde{Dist}_{C1-A6} * \tilde{Dist}_{C1-B1}} \right)$$

Idiosyncrasy #1: Fuzzy Lines

- Given a TFN, that corresponds to a line, what are its endpoints?



Idiosyncrasy #2: Fuzzy Vertices

- When two or more fuzzy lines intersect, what is the shape of the resulting fuzzy vertice?



Example: CuInSe_2

ABC_2 compound

- 13 A Atoms
- 10 B Atoms
- 8 C Atoms

QuickTime™ and a
Animation decompressor
are needed to see this picture.

QuickTime™ and a
Animation decompressor
are needed to see this picture.



We isolate a single Se
atom for our analysis.

CuInSe₂: Fuzzy Bond Lengths

Bond	Fuzzy (Å)	Knight, 1992 (Å)	μ
Se-Cu	{2.390, 2.442, 2.527}	2.4337	0.836
Se-Cu	{2.394, 2.442, 2.518}	2.4337	0.821
Se-In	{2.488, 2.574, 2.632}	2.5893	0.730
Se-In	{2.494, 2.574, 2.624}	2.5893	0.689

CuInSe₂: Fuzzy Bond Angles

Angle	Fuzzy (°)	Knight (°)	μ
Cu-Se-Cu	{109.481, 113.932, 115.947}	114.870	0.534
Cu-Se-In	{106.234, 109.144, 111.512}	108.902	0.917
Cu-Se-In	{105.924, 109.426, 112.538}	109.462	0.988
Cu-Se-In	{105.958, 109.426, 112.582}	109.462	0.989
Cu-Se-In	{106.608, 109.144, 111.430}	108.902	0.905
In-Se-In	{102.705, 105.423, 109.474}	104.760	0.756

Example: YBCO

Still In Progress!!

Conclusions

- Fuzzy logic is an innovative method for representing imprecision.
- When applied to molecular modeling, it can represent measurement error, structural defects, and thermal and vibrational variations.

Imaging Studies and Density Functional Analysis of Surfaces and Interfaces: Comparison of Theory and Experiment

John F. Maguire and Steven R. Le Clair

Air Force Research Laboratory,
Materials and Manufacturing Directorate,
Wright-Patterson AFB , Ohio.

ABSTRACT

The ability to see an image allows the human brain to engage in a range of reasoning and cognitive functions which exceed that which is possible with purely numerical data. Similarly, the ability to assess the solution to a coupled set of mathematical equations lies beyond symbolic manipulation and approaches creativity in advanced physical reasoning. The ability to develop imaging technology over an appropriate range of energy and momentum transfer coupled with an ability to reason about and interpret such images is critical to the development of A.I. for process control. In this paper, we present the results of two imaging studies which have been directed towards developing surface and near surface imaging technology for material and process control.

The first was an experimental model system chosen to simulate the near surface region of a polymeric composite. Polymers were modeled as an homologous series of oligomeric alkane fragments which were deposited on the [0001] plane (basal) of highly-ordered pyrolytic graphite. Raman and scanning tunneling spectroscopy reveal that matter in the surface and near surface region takes on a qualitatively different structure from the bulk material. The alkanes form a self-assembled monolayer on the surface and adopt a characteristic rank and file structure on the surface.

In the second system, Raman imaging and STM imaging studies were applied to the surface of thin superconducting films of Yttrium Barium Copper Oxide which were deposited on the [111] face of a lanthanum oxide substrate using pulsed laser deposition. There is evidence that the performance of these films is influenced by the stoichiometric ratios and it is therefore of interest to determine if the surface material is homogeneous or whether there is evidence of speciation and/or aggregation. We have observed Raman images which are consistent with macroscopic structure formation over micron distance scales.

The underlying mechanisms for structure formation are discussed in terms of the symmetry of local fluctuations and the symmetry breaking role of the interface. Model results are illustrated by application of density functional theory to the surface and near surface region.

Imaging Studies and Density Functional Analysis of Surfaces and Interfaces

Comparison of Theory and Experiment

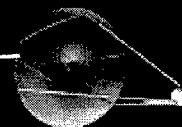
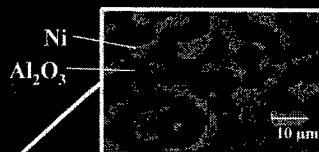
John F. Maguire & Steven R. LeClair
Air Force Research Laboratory

Material

Examples - Materials & Needs

Co-continuous Composites for Thermal Protection
Metal-Ceramic, Metal-Polymer, Ceramic-Polymer Composites
and Light-weight Radiation Resistant Materials

Optical Coatings for Deformable Mirrors
Reaction Diffusion Coatings for In-service Reconstruction
Adaptive Coatings for Air-Space Missions
Films for Light-weight Systems for Directed Energy



YBCO tapes

Zinc Oxide
Wear Coatings

The challenge - to accelerate space materials research -
via prototyping of new materials/processes
based upon mappings of process designs to controlled results.

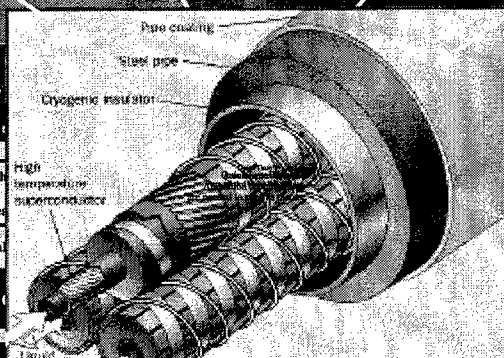
Overview

MOLECULES < MESOMATERIALS < NANOMATERIALS

Tape to Wire
Requires Spiral Tw
1-2 rev/3" length
With 0.2% film strain

1 tape - 2 μm th
or
2 tapes - 1 μm ea

Metal
No



should be 10:1
(O to insulator)

YBCO - HTS
CeO₂ - Buffer
Invar Alloy - Substrate

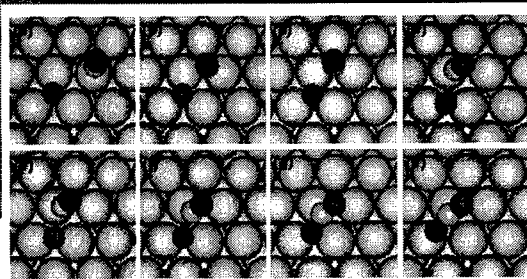
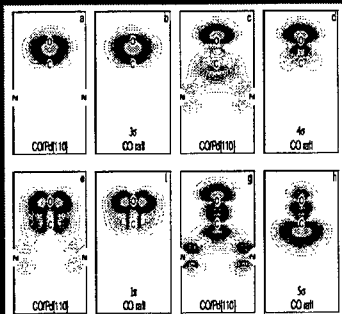
Time dependent Schrödinger equation solved for N body system

**Molecular Mechanics - PLD growth of
CeO₂ ISLANDS OF KNOWLEDGE**

Quantum Mechanical Calculations for Large Systems

ab initio calculations of
catalytic reactions
(e.g., Car-Parrinello Code)

CO to CO₂
Orbital Mixing Model



Physical Model Sequence of
CO to CO₂ Chemisorption on Pd

Challenge: enable the approximation
of these calculations for use in
near real-time process simulation

Ab initio vs Hybrid Approach

Size of Problem	# of Atoms	# of Forces (Inter-nuclear)	Time Cray	Time PC
Useful to study Quantum Effect	15	105	300 hours 0.6 fs/step, ~1000 steps	-
Cray doing Process Simulation	2×10^5	2×10^{10}	6×10^{10} hours	-
PC doing Process Simulation	2×10^5	30 local forces	-	0.5 hours local vs global interactions

A. Alavi et al. Phys. Rev. Lett. 80, 16 p. 3650, 1998

10^{12} advantage

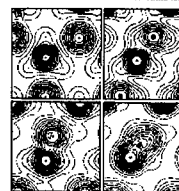


FIG. 3. Electron density isosurfaces parallel to the surface. The plane of the isosurface is chosen as midway between the z intercepts of the C atom and H_{10} . The contours are logarithmically spaced. (a) Black contours show the highest electron density (H_{10}). The isosurfaces spacing (1) is even better. Densities are in atomic units ($a.u.$). (b) CO_2 ($2 \times 2 \times 2$ a.u.) (c) CO_2 ($2 \times 2 \times 2$ a.u.) (d) CO_2 ($2 \times 2 \times 2$ a.u.)

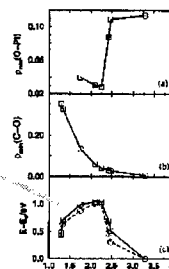


FIG. 4. (a) and (b) show the change in the overlap of the atomic orbitals (AOs) along the reaction pathway. $\rho_{A,B}(A,B)$ is the lower value of the electron density along the axis from atom A to atom B . The change in $\rho_{A,B}$ is a convenient measure of the change in the strength of the chemical bond between A and B . The electron densities in (a) and (b) are shown. (c) The energies of the eight configurations shown in Fig. 2 relative to the lowest energy configuration. Model line: LDA; solid line: GGA.

'Hybrid' Molecular-Process Simulation via

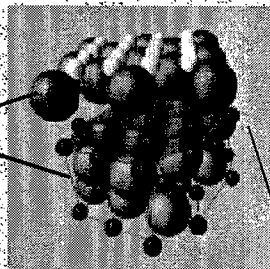
the coupling of ab initio calculations and cellular automata

to model the local forces and to parallelize and distribute computation

Speed - to bridge the gap, we need 90% of benefit in 1% of the time

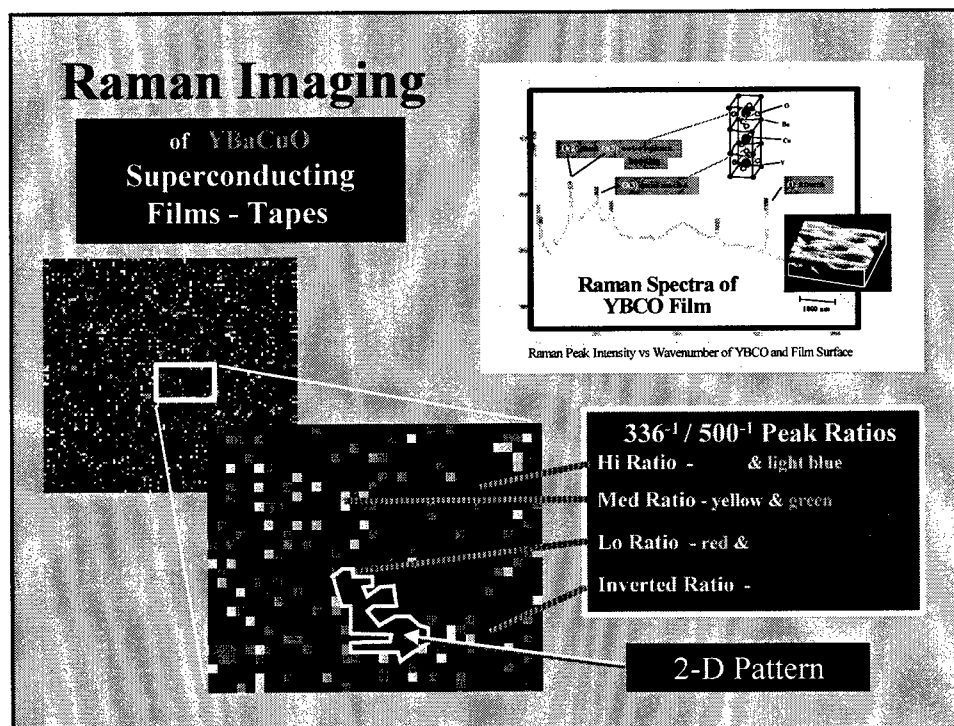
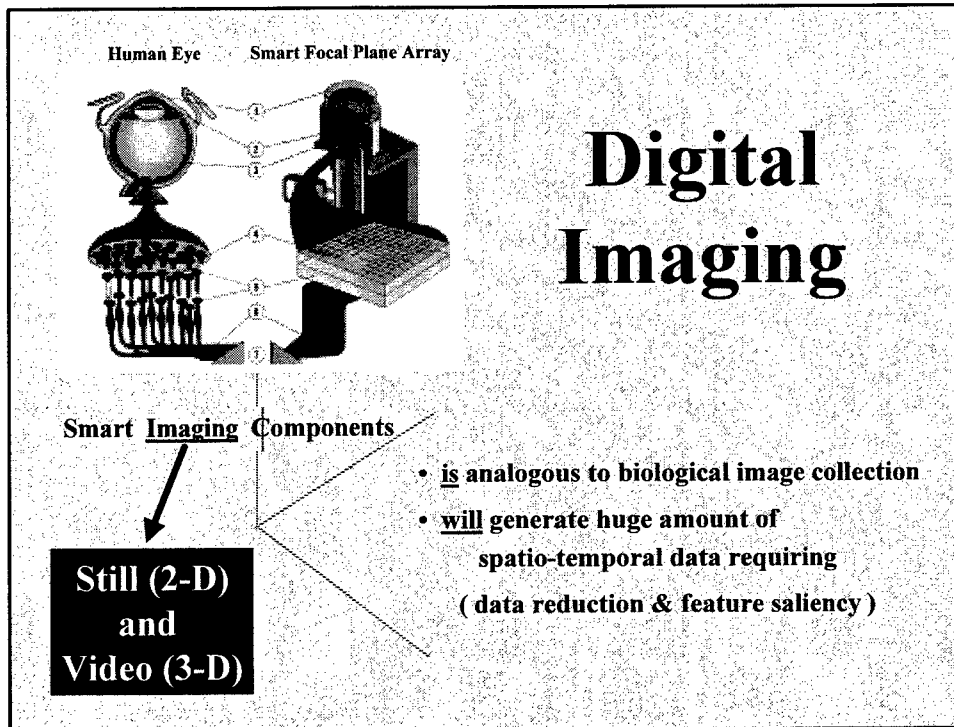
th process - to y, composition less

Material Interface Design:
dissimilar materials, varying structure, electro-negativity, ...



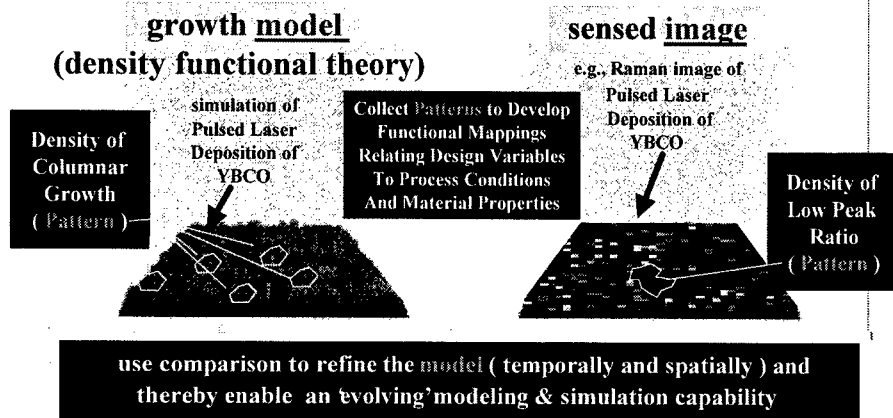
« 1000 x 1000 atoms subs

Accuracy - needs to be sufficient for process control, but essential for production of long length (YBCO tapes)



'Evolve' Design-Control Mappings

via mapping of in situ *sensor imaging* (Raman) picture/video
and compare against model-generated picture/video



* Virtual implies the ability to interpolate/extrapolate the sensed image from a yet-to-be-made material system

Modeling Gas By-Products from MO-CVD Thin-Film Depositions

J.G. Jones and P.D. Jero

Air Force Research Laboratory, Materials Directorate
Wright-Patterson AFB, OH, 45433-7746

This work seeks to develop an MO-CVD system using *in situ* sensors and automated process control to deposit controlled, reproducible oxide fiber coatings (e.g. $\text{LaAl}_{11}\text{O}_{18}$). A CVD system capable of continuous fiber coating has been assembled which employs liquid precursor delivery for precise precursor stoichiometry control and inert gas seals for near atmospheric operation. *In-situ* thermocouples and a mass spectrometer are used for process measurements. All of the system parameters are logged by and controllable by computer. A fuzzy logic controller was developed to control the O_2 flow rate based on the desired temperature or gas composition. A neural model of the process will be presented which translates the process settings into expected gas compositions.

Keywords: Process control, Fuzzy logic, Modeling, Neural Networks, Process identification

Modeling Gas Byproducts From MOCVD Thin-Film Depositions

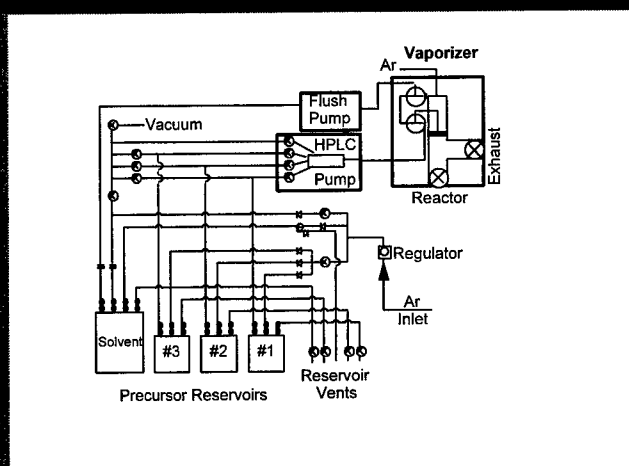
J. G. Jones and P.D. Jero

Air Force Research Laboratory
Materials & Manufacturing Directorate

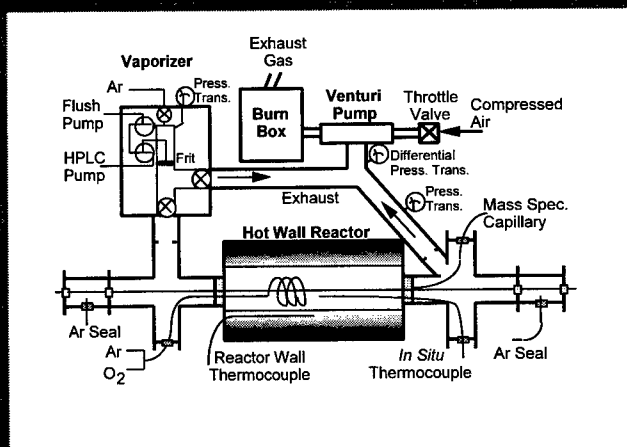
MOCVD Fiber Coating Apparatus



Liquid Precursor Delivery System

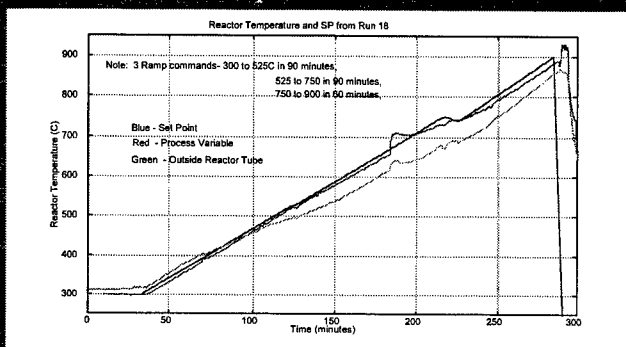


MOCVD Hot Wall Reactor System

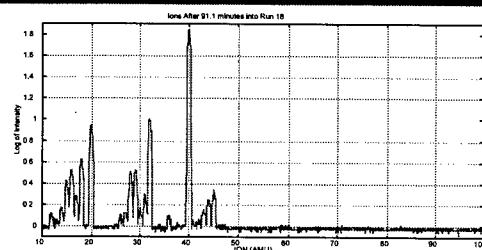


Ramp of Reactor Temperature

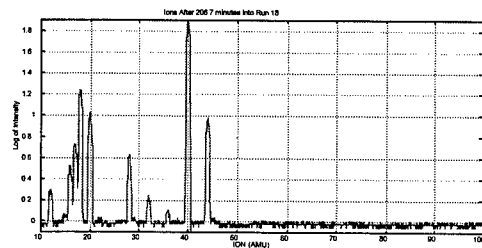
External Wall Temperature Control &
In Situ Temperature Measurement



Gas Phase Products

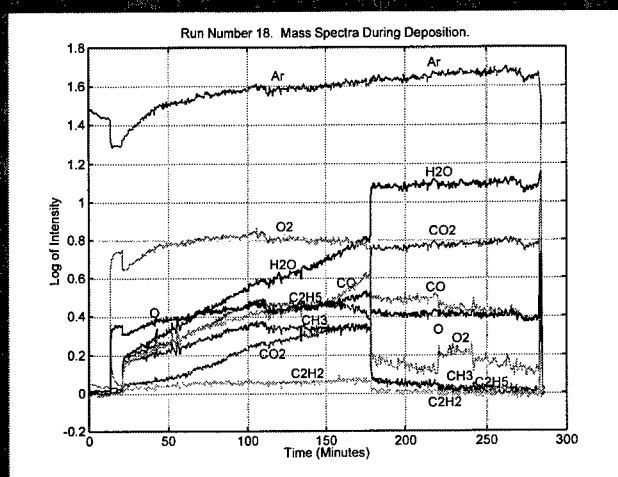


~450 C :
incomplete
combustion

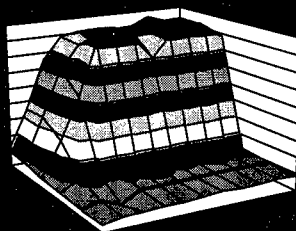


~710 C :
complete
combustion

In Situ Gas Phase Products



Intensity of H₂O in Gas Phase



Imaging for Process Optimization and Control

(abstracts and viewgraphs)

4

Nondestructive Imaging of Surface and Sub-Surface Defects in Thin-Films with Super Spatial Resolution Using Evanescent Microwave Fields

Massood Tabib-Azar

Case Western Reserve University
Cleveland, Ohio 44106
Tel: 216-368-6431

With the current spatial resolution of 4 μm at 1 GHz or 0.4 μm at 10 GHz, the evanescent microwave probe (EMP) is capable of mapping non-uniformities and defects in a variety of materials including insulators, semiconductors, metals, biological, and botanical samples. Due to its large center frequency of operation, EMP can be operated with very fast scan rates approaching 10 cm/s. Since the method does not require physical contact with the sample, it can also be used to study "sticky", as well as, hot and cold surfaces. According to our recent studies, the spatial resolution of EMP can be improved by a factor of 40 in certain materials. Its principle of operation is based on producing decaying electromagnetic fields near a discontinuity in a micro-strip-line waveguide. EMP is a two-dimensional planar structure and it is possible to design and fabricate many parallel EMPs operating at different frequencies on silicon cantilever beams. Such an integrated parallel probe assembly will enable efficient, fast and hyper-spectral mapping of moving samples. There are manufacturing problems that can be addressed and solved by using fast scanning probes in an assembly or manufacturing environment to provide real-time information and feedback regarding the quality of deposited films. In this paper, we will discuss some recent experimental results and probe characteristics.

Nondestructive Imaging of Surface and Subsurface Defects in Thin-Films with Super Spatial Resolution using Evanescent Microwave Fields

M. Tabib-Azar* and S. R. LeClair**

*Electrical Engineering and Computer Science Dept. &
Macromolecular & Physics Depts.
Case Western Reserve University
Cleveland, OH 44106 (e-mail: mxt7@po.cwru.edu)

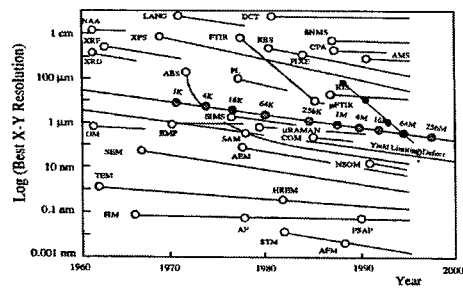
****Air Force Research Laboratory
Materials & Manufacturing Directorate
2977 P. Street, Suite 13, Wright-Patterson AFB, OH**

Supported by grants from: Wright-Patterson Air Force Base, CAMP, and NASA

Outline

- I. Evanescent Sensing Methods?
- II. Resolution
- III. Results
- IV. Future Work

History of Analytical Techniques in Semiconductor Manufacturing



Super-Resolution Non-Destructive Imaging of Defects, and Non-uniformities in Metals, Composites, Semiconductors, and Dielectrics Using Evanescent Microwave Probes

Massood Tabib-Azar

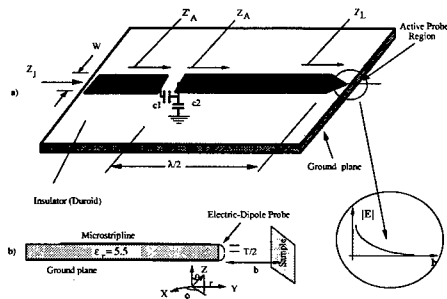
**Electrical, System, Computer Engineering and Science Department
& Macromolecular Science
& Physics Departments**

Case Western Reserve University,
Cleveland, OH 44106

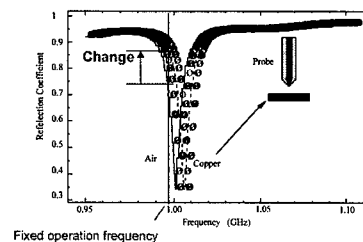
(phone: 216-368-6431, fax: x6039, e-mail: mxt7@po.cwru.edu)

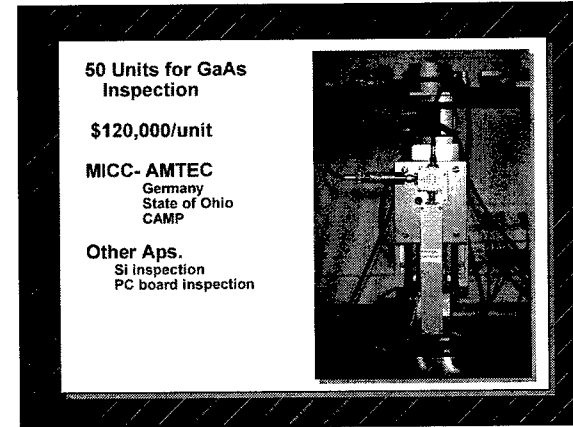
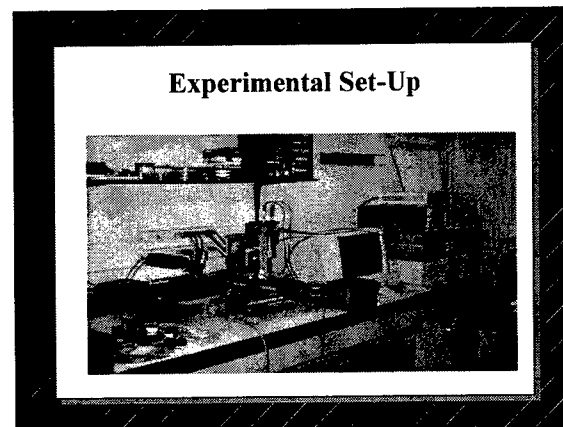
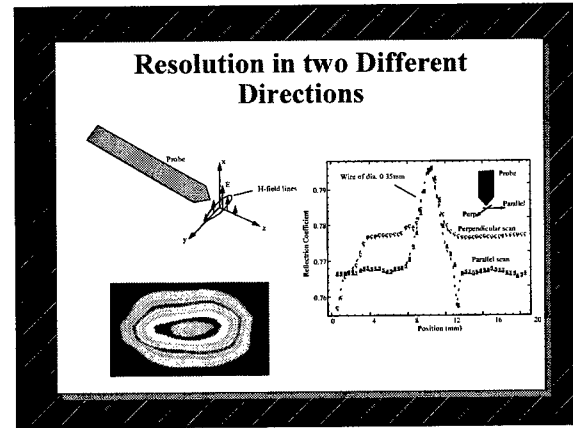
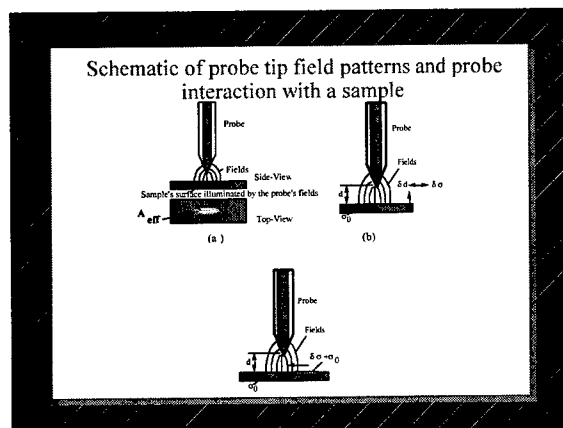
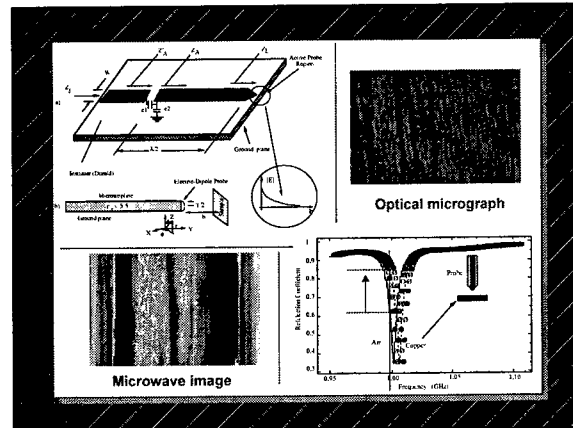
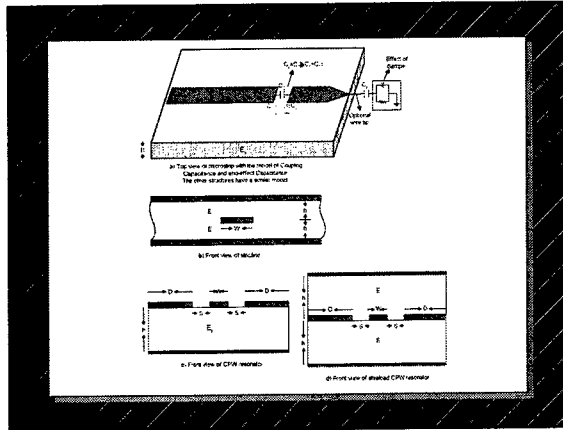
Supported by grants from: Wright-Patterson Air Force Base

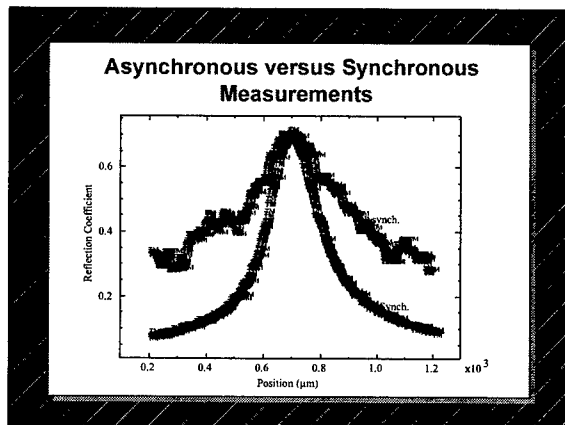
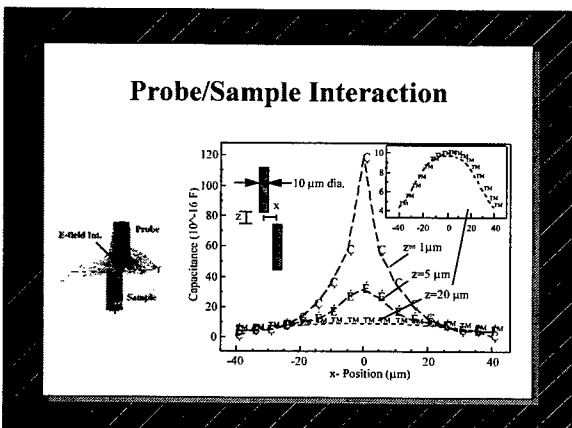
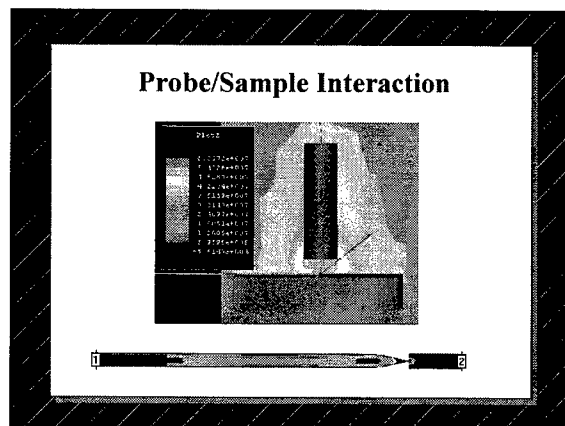
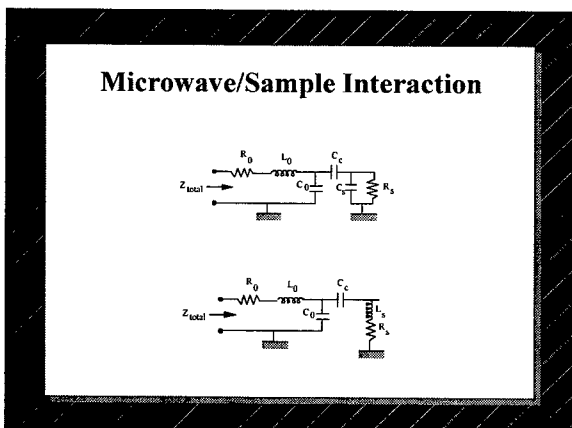
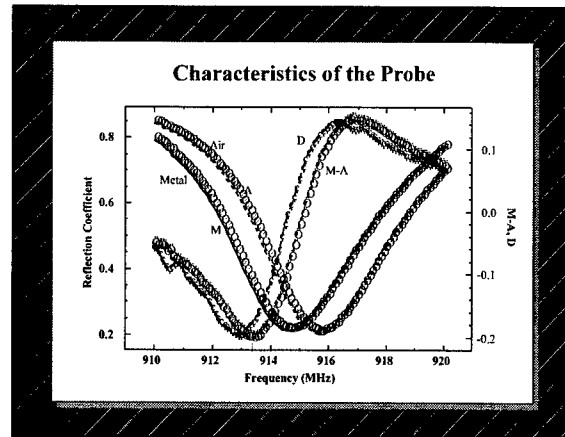
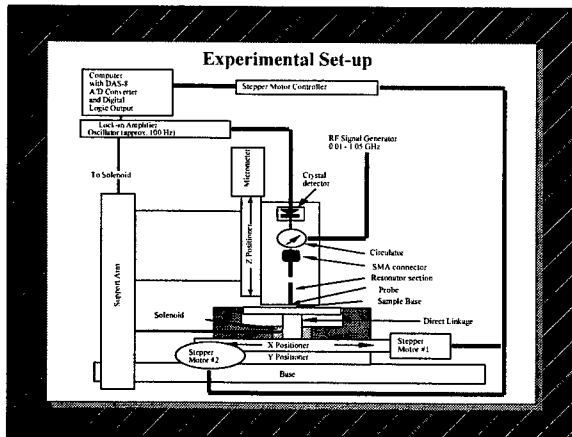
Description of the Method



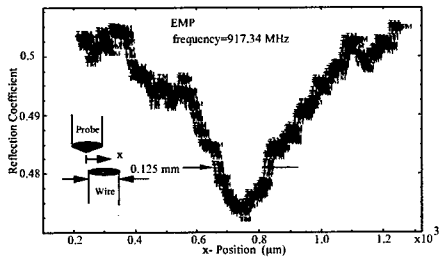
Characteristics of the Probe



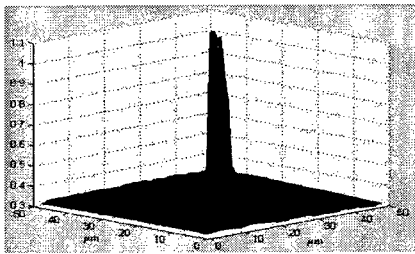
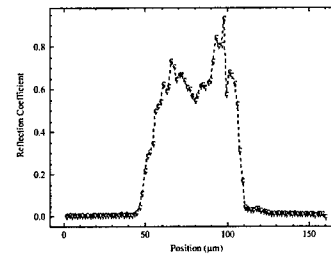




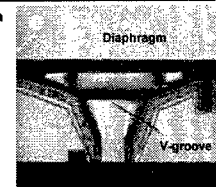
A linescan over top surface of a 0.125 mm dia. Wire



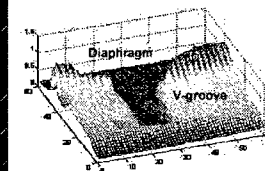
Features at the tip of a 50 μm wire



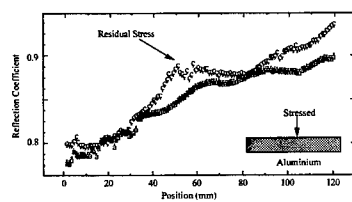
Optical micrograph



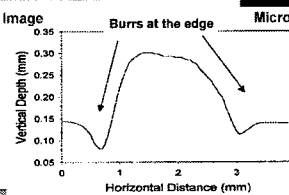
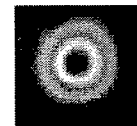
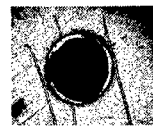
Microwave Image



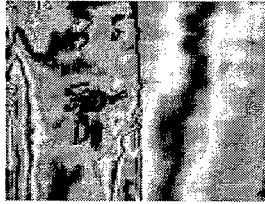
Detection of Residual Stress



Topographical Map of a Hole in Brass



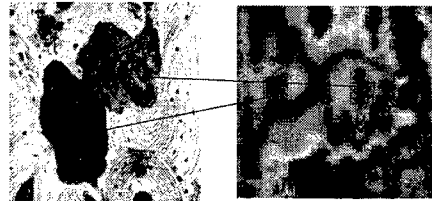
Evanescent Microwave Map of Delamination in Carbon Composites



Delaminated

Intact

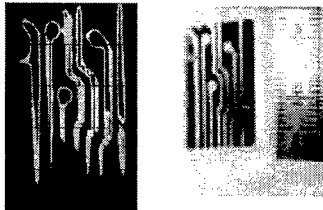
Microwave Images of Bone



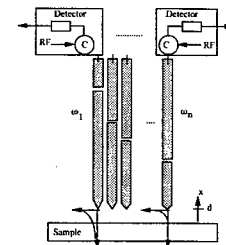
Acoustic Image

Evanescent Microwave Image

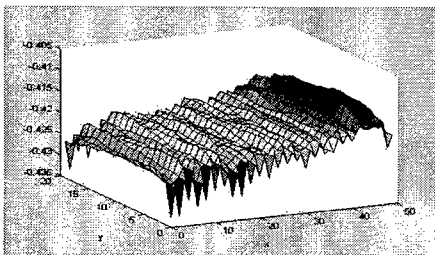
Printed Circuit Board



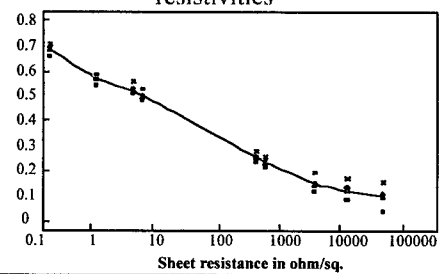
Probes operating at different frequencies will be used to obtain information regarding the distribution of impurities and other parameters as a function of depth.



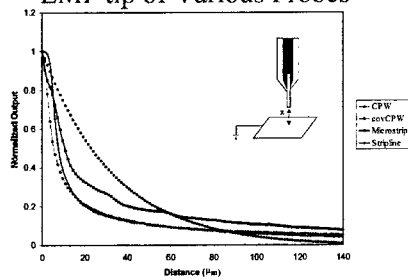
Hyperspectral image of PLD YBCO extracted from EMP scans at 1 and 10 GHz



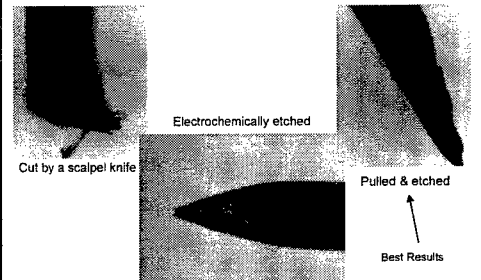
Evanescent microwave probe calibration using silicon samples with different resistivities



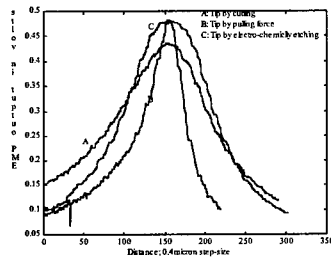
Characteristic Decay Lengths Near the EMP tip of Various Probes



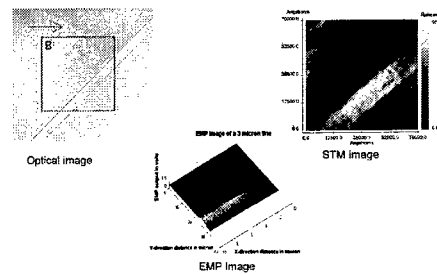
Wire-Tip Preparation



Wire-Tip Effects

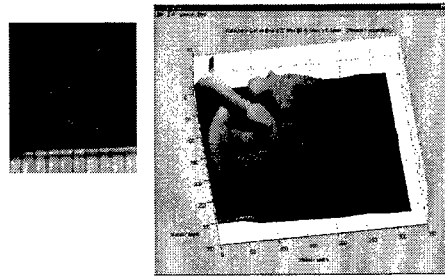


Optical, STM, and EMP Images of a 3µm Line



EMP scans of ethylene propylene rubbers with rough and smooth surfaces

Photograph of a metallic pattern (dia. of 7 mm) on a PC board and its EMP image



Investigation of Raman Imaging for Advanced Control of YbCO Cool-Down Processing Using Pulsed Laser Deposition

J.D. Busbee ^{*1}, R.R. Biggers ^{*}, J.G. Jones ^{*}, D.V. Dempsey ^{*2}, G. Kozlowski ^{**}

^{*} Air Force Research Laboratory, Materials Directorate
Wright-Patterson AFB, OH, 45433-7746, USA

^{**} Air Force Research Laboratory, PRP, Wright-Patterson AFB, OH, USA

¹ Technical Management Concepts, Inc., Beavercreek, OH 45434

² University of Dayton Research Institute, Dayton, OH 45409

Pulsed Laser Deposition (PLD) is a versatile, complex thin film deposition process that has been shown to be capable of creating high quality YbCO films. However, the promise of the technique has been hampered by a lack of process understanding and visibility into in-situ processes. In the past, attempts at controlling the PLD process have not been highly successful due to a lack of direct feedback from the surface of the film. This paper investigates using Raman spectroscopy to provide feedback from the evolving film during cool-down after deposition. In-situ spectra taken under controlled conditions are examined to understand the effect of environmental parameters on film properties as well to provide insight into evolution of the film microstructure during cool-down.

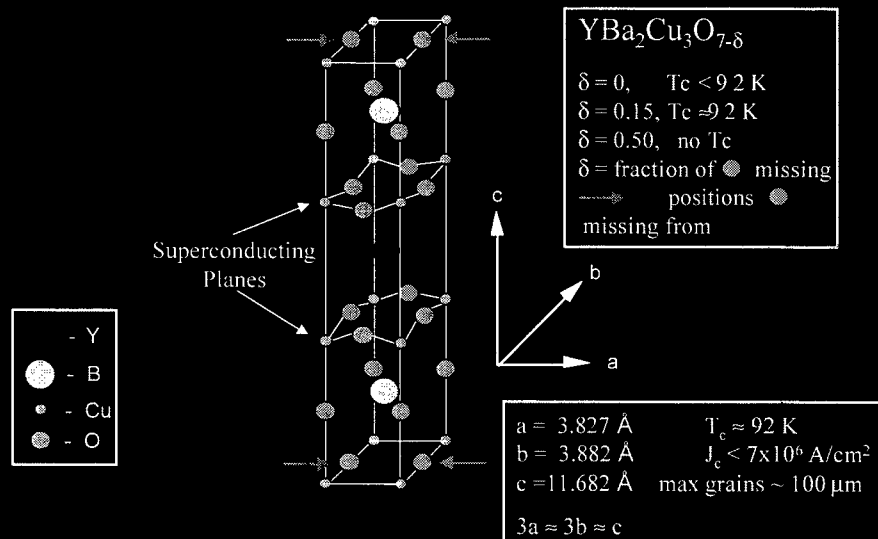
Keywords: Process control, PLD, Raman Spectroscopy, Superconductors, YBCO

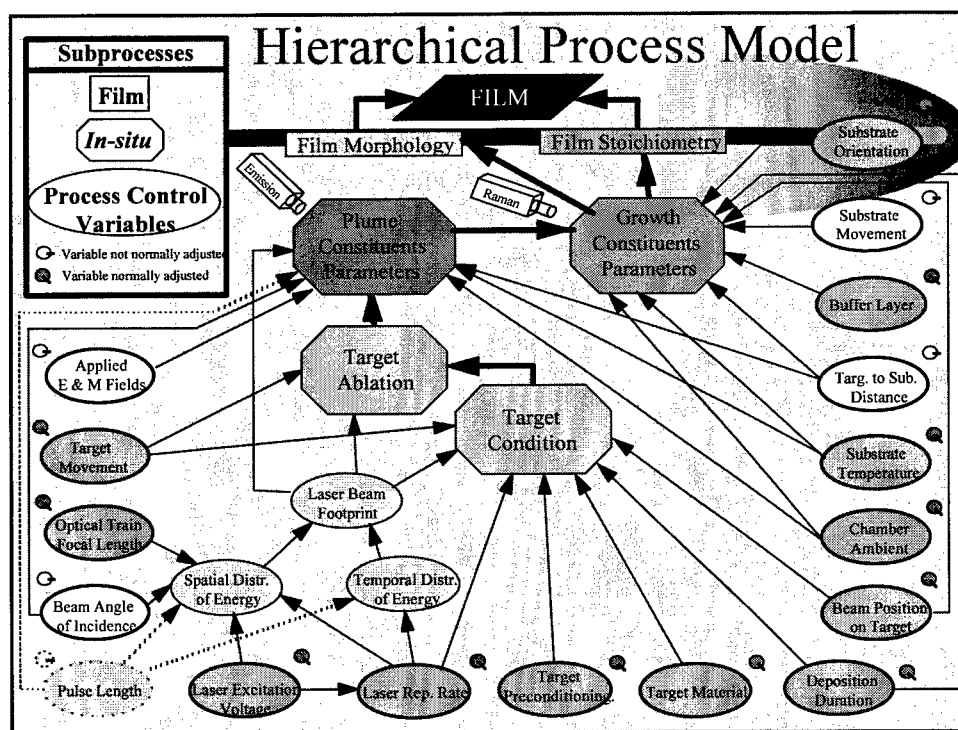
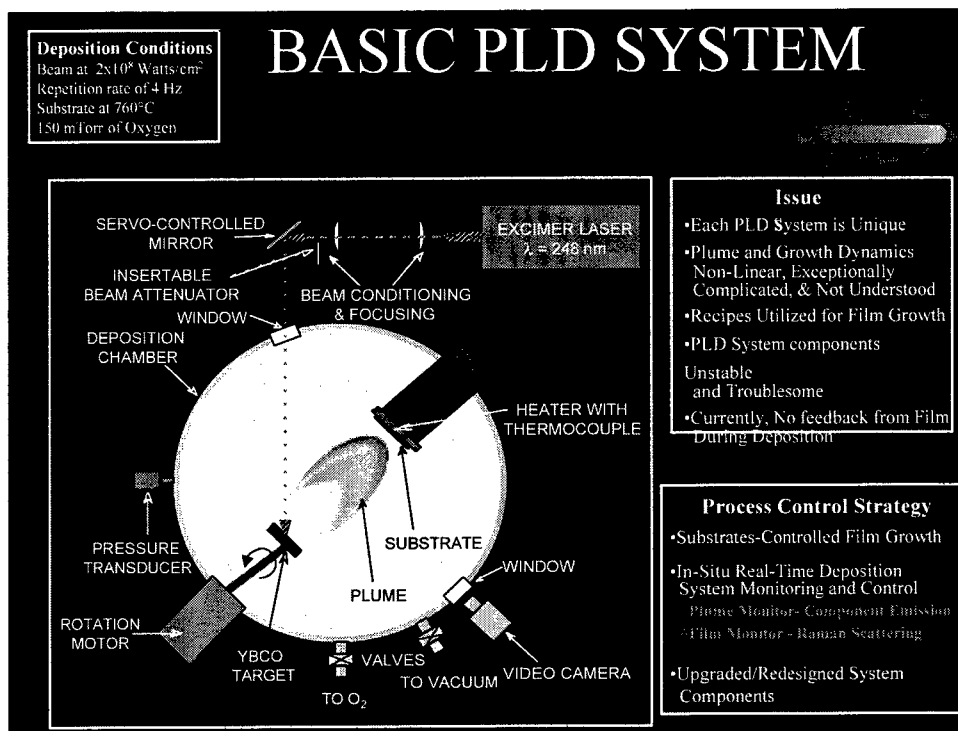
Investigation of Raman Imaging for Advanced Control of YBCO CoolDown Processing using Pulsed Laser Deposition

John Busbee
Air Force Research Lab
Materials Directorate
AFRL/MLMR

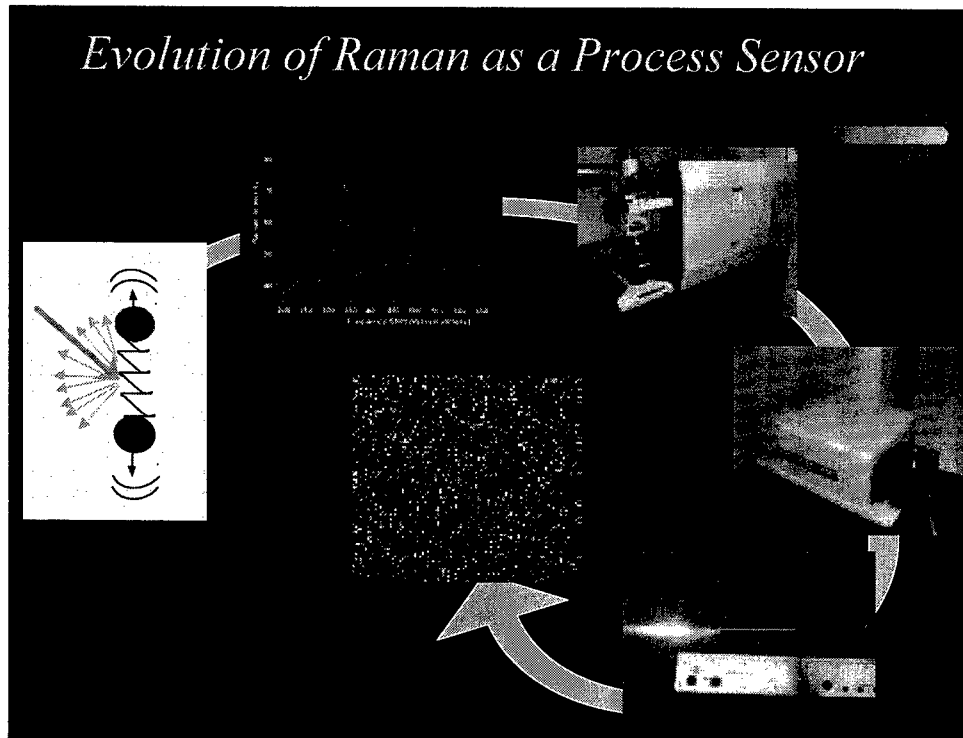
MATERIAL UNDER DEVELOPMENT

YBa₂Cu₃O₇ ORTHORHOMBIC UNIT CELL

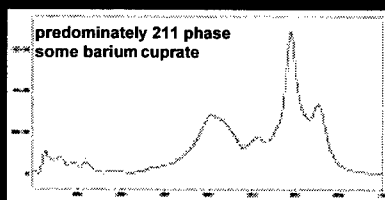




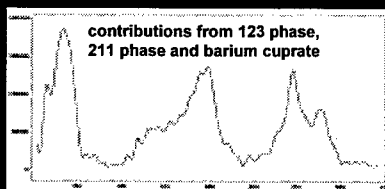
Evolution of Raman as a Process Sensor



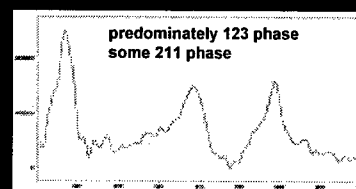
Stoichiometry of $YBa_2Cu_3O_{7-x}$



$T_c = 80K$

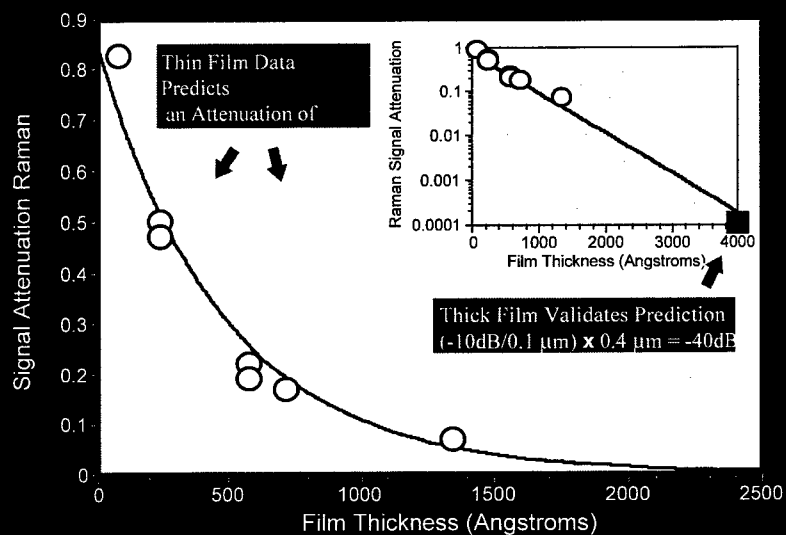


$T_c = 87K$

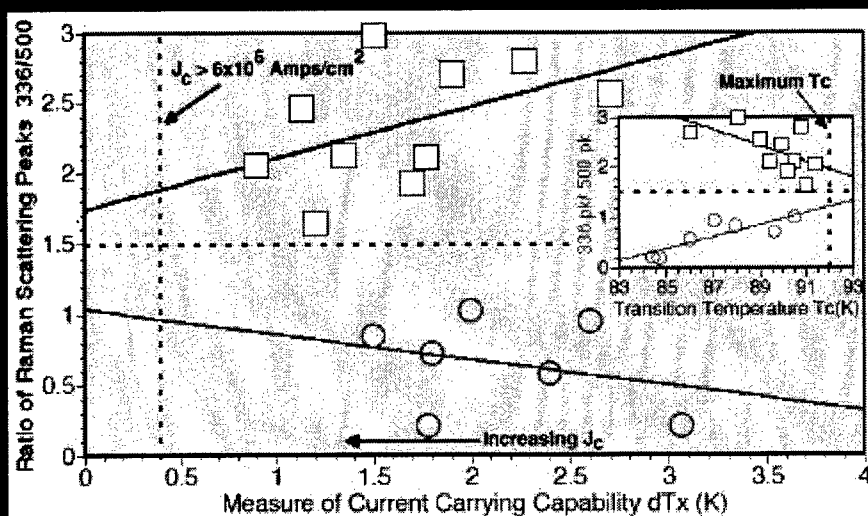


$T_c = 91K$

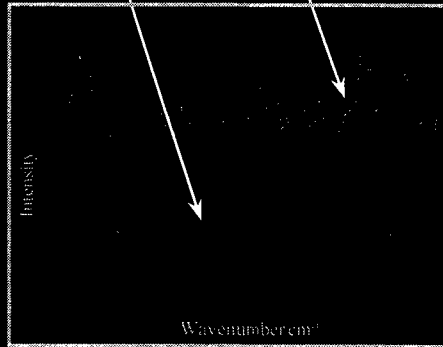
YBCO THICKNESS MEASUREMENTS USING RAMAN SPECTROSCOPY



CORRELATION OF RAMAN PEAK RATIOS WITH CRITICAL CURRENT DENSITY, J_c , AND TRANSITION TEMPERATURE T_c

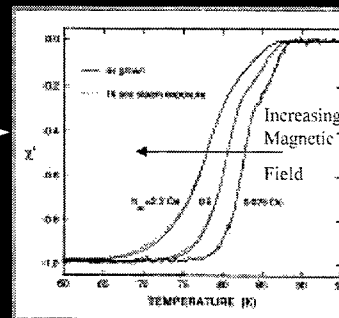


Raman Spectra of YBCO Thin Film Before and After Exposure to Steam



Raman Spectra Show Significant
Surface Degradation

AC Susceptibility Tests Show Film Qualities
(T_c , J_c) Virtually Unchanged



Results:

Current Carried Primarily in Deeper Layers of YBCO
&
Raman Best Used as *In Situ* Sensor

Process Control via Gaze Detection Technology

Jaihie Kim*, Gang Ryung Park* and Steven LeClair**

* Dept. of Electronic Engineering,

Yonsei University, Seoul 120-749, Korea

Tel: 82-2-361-2869 Fax: 82-2-362-0413 Email: jhkim@bubble.yonsei.ac.kr

** Materials Process Design Branch, Manufacturing Technology Division

Materials & Manufacturing Directorate, Air Force Research Laboratory, U.S.A.

Tel: 1-937-255-8787

Gaze detection technology (GDT) is considered to be of potential benefit to automate certain materials processing research tasks and enable more remote control of manufacturing processes. Gaze detection is defined as the ability to assess a human user's direction of view and/or a position focus on a computer monitor screen via computer vision techniques. The research reported herein will evaluate the suitability of gaze detection technology in the remote control of a computer via a monitor screen and a gaze detection system. Monitoring of the user's face movement is achieved by using a camera included in the system. A user-interface depicting a process control environment is artificially designed for this work and results show that GDT is very useful in automating and/or enabling remote actuation of various process management tasks.

To be useful, GDT must enable remote control of three tasks. First, it must locate the position, on a monitor screen, where the user is looking, i.e., place a cursor at that position. Second, it must detect in which direction a user's face is moving, i.e., to drag the cursor in that direction or to stop the cursor when it is moving in the reverse direction. These two tasks may be invoked sequentially or independently. When they are invoked sequentially, the latter usually compensates for the former positioning error. It should be noted that these two tasks may be independently useful for various purposes. Finally, as the camera is tracking one of the eyes, e.g., the left eye of the user, and the user winks, then the GDT system is capable of actuating some device. Thus, the GDT system is capable of placing an indicator at a specific position on a screen, dragging and clicking on an icon, etc., remotely without using any direct contact.

GDT will be discussed in the context of materials processing, and, more specifically, to address three aspects of process discovery. First, it enables gaze control (e.g., allowing a user to visually open and close windows) of the monitor panel whether or not the user's hands are busy doing other things. This type of control is enabled by the camera monitoring the human gazing at the monitor and detecting the winking of the left eye. Second, the system determines if the user is not present or within distance of an input device, and if not, then it displays and invokes a special function. For example, gaze control could activate a process discovery program to perform parameter assessment to detect either an anomalous condition or parameter perturbation of interest, and records the time of occurrence and generate an audible or visual alarm. The user must then decide, using GAZE CONTROL, to replay the screen or window display to observe the detected anomaly or perturbation.

A GDT system will be discussed which is capable of determining if a user is focusing on a particular spot of interest on the monitor. If not, the GDT system will attract the user's attention to that spot (again, coupling GAZE CONTROL to a program for generating an audible or visual signal) and direct the user's attention to the spot of interest. An example use of this capability might be to GAZE CONTROL the temporal display of a particular parameter which contains current parameter values only, or highlight particular variables associated with a detected fault or unusual process condition.

Process Control via Gaze Detection Technology

1999. 7

Jaihie Kim*, Kang Ryoung Park*, Steven LeClair**

*Dept. of Electrical and Computer Engineering, Yonsei University
Seoul 120-749, Korea, E-mail : jhkim@bubble.yonsei.ac.kr
**Materials Process Design Branch, Manufacturing Technology Division
Materials & Manufacturing Directorate Air Force Research Laboratory, U.S.A.

Outline

- I What is Gaze Detection Technology?
- I Implementation Environments
- I Gaze Detection by Face and Eye Movements
- I Gaze Position Detection Techniques
- I Additional Techniques for Computer Interface
- I Protocols for Computer Interface
- I Application of Gaze Detection Technology to Process Controls
- I Other Applications
- I Future Works

Dept. of Electrical and Computer Engineering

Outline

What is Gaze Detection Technology?



Dept. of Electrical and Computer Engineering

What is Gaze Detection Technology?



Dept. of Electrical and Computer Engineering

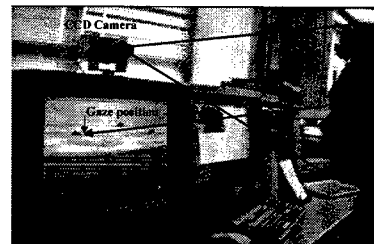
What is Gaze Detection Technology?



Dept. of Electrical and Computer Engineering

What is Gaze Detection Technology?

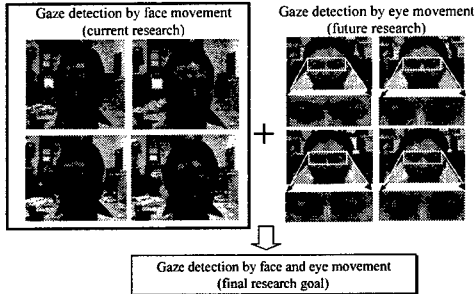
Implementation Environment



Dept. of Electrical and Computer Engineering

Implementation Environment

Gaze Detection by Face and Eye Movements

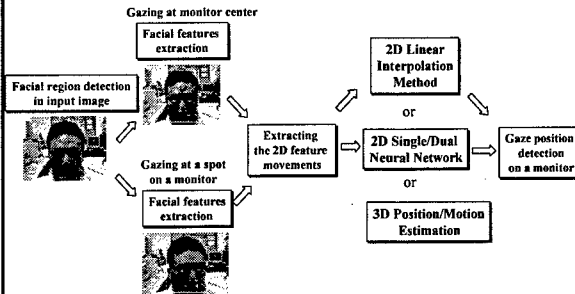


Dept. of Electrical and Computer Engineering

Gaze Detection by Face and Eye Movement

Gaze Position Detection Techniques

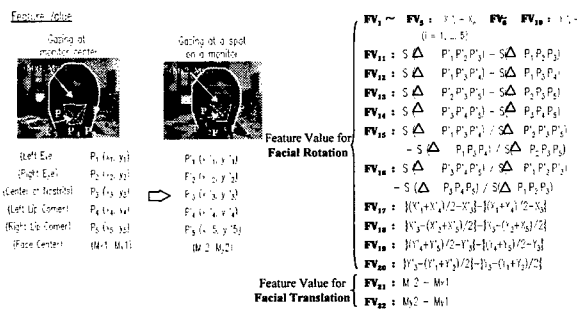
1 Overview of our gaze detection techniques



Dept. of Electrical and Computer Engineering

Gaze Position Detection Techniques

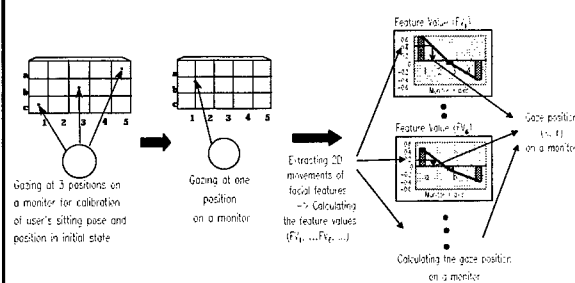
Feature Value for Gaze Position Detection



Dept. of Electrical and Computer Engineering

Feature Value for Gaze Position Detection

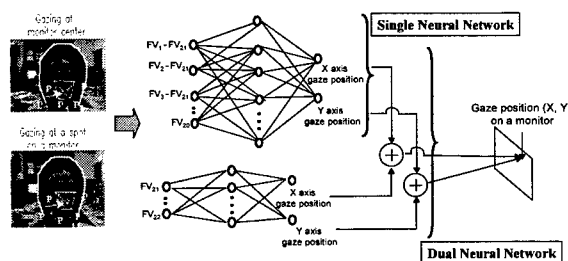
Linear Interpolation Method



Dept. of Electrical and Computer Engineering

Linear Interpolation Method

2D Single/Dual Neural Network - (1)



Dept. of Electrical and Computer Engineering

2D Single/Dual Neural Network

2D Single/Dual Neural Network - (2)

1 Normalization of Feature Value

To **normalize the distances of users from the monitor**, a user is asked to look at the center, a right-most upper point and a left-most lower point predefined on the monitor when he begins to use the system.

We obtain maximum and minimum feature values from this, and each feature values are normalized by these maximum and minimum feature values as follows.

$$\overline{FV}_i = \frac{FV_i}{\text{Max}(FV_i) - \text{Min}(FV_i)} \quad (i=1 \sim 10)$$

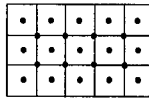
Dept. of Electrical and Computer Engineering

2D Single/Dual Neural Network

2D Single/Dual Neural Network - (3)

1 Training data for Neural Network

Training data were obtained from 10 users looking at 23 known points on a 19" monitor



- Gaze positions for training

1 Generalized Delta Rule for Neural Network Training

$$W_{kj}(t+1) = W_{kj}(t) + \eta \delta_{pk} i_{pj}$$

$W_{kj}(t)$: the weight values of each node when the iteration number is t

η : learning rate parameter ($=0.01$)

δ_{pk} : the error term in hidden units or output units

i_{pj} : the calculated output in input layer or hidden layer

Dept. of Electrical and Computer Engineering

2D Single/Dual Neural Network

Comparisons & Experimental Results

Gaze position errors are calculated by the differences between the real gaze position and the calculated position

(inches)

Method	Linear Interpolation	Single Neural Net	Dual Neural Net
RMS error	1.84	.64	1.7

? Gaze detection errors for experimental data including facial rotation
(5 persons * 22 gaze positions = 110 test data)

Method	Linear Interpolation	Single Neural Net	Dual Neural Net
RMS error	4.54	4.40	3.4

? Gaze detection errors for experimental data including facial rotation and translation
(5 persons * 22 gaze positions = 110 test data)

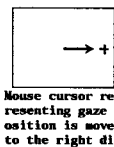
Dept. of Electrical and Computer Engineering

Comparison & Experimental Results

Additional Techniques for Computer Interface - (1)

1 Mouse Dragging by Facial Movement

In order to compensate the gaze detection error, the user moves his face to the desired point on the monitor, then the monitor cursor placed by gazing is moved toward the point.



The facial movements are measured by the 2D movements of facial features when a user rotates and translates his face

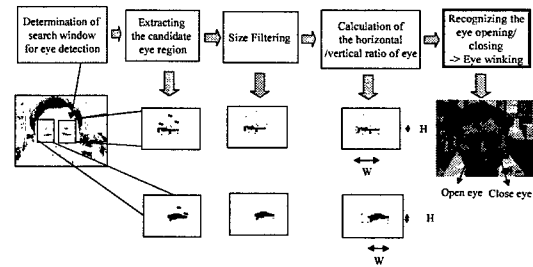
The current facial feature positions are predicted from previous feature positions, in order to reduce processing time.

Dept. of Electrical and Computer Engineering

Additional Techniques for Computer Interface

Additional Techniques for Computer Interface - (2)

1 Clicking a Mouse Button by Eye Winking

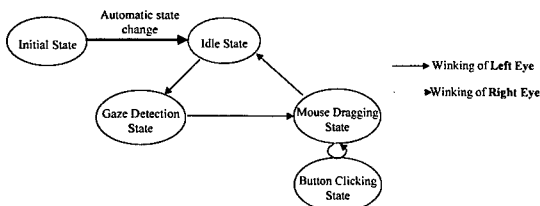


Dept. of Electrical and Computer Engineering

Additional Techniques for Computer Interface

Protocols for Computer Interface - (1)

Case 1: Interface for custom-made S/W



Initial State: The state when system is starting

Idle State: The state when a user does not want to use gaze detection system

Gaze Detection State: The state when the user gazes at a spot on a monitor

Mouse Dragging State: The state when the user drags the mouse cursor by facial movement

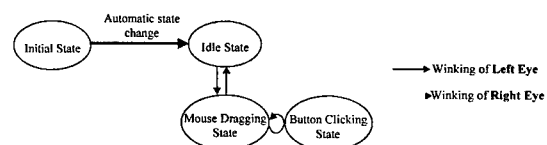
Button Clicking State: The state when the user commands to the current mouse cursor

Dept. of Electrical and Computer Engineering

Protocols for Computer Interface

Protocols for Computer Interface - (2)

Case 2: Interface for ready-made S/W



Initial State: The state when system is starting

Idle State: The state when a user wants to use mechanical mouse

Mouse Dragging State: The state when the user drags the mouse cursor by facial movement

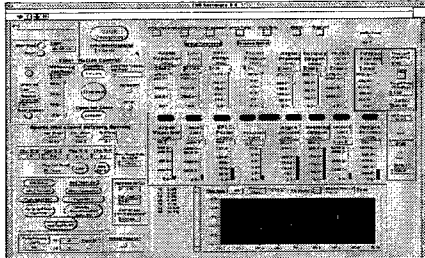
Button Clicking State: The state when the user commands to the current mouse cursor

Dept. of Electrical and Computer Engineering

Protocols for Computer Interface

Application of Gaze Detection Technology to Process Controls - (1)

1 A user interface of custom-made CVD(Chemical Vaporized Deposition) program

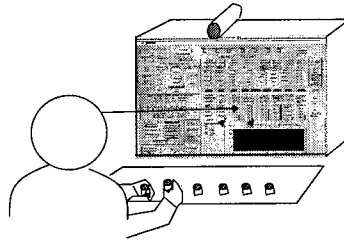


Dept. of Electrical and Computer Engineering

Application of Gaze Detection Technology to Process Controls

Application of Gaze Detection Technology to Process Controls - (2)

1 In its first function, gaze detection technology enables a user to control another process, even when both hands are busy in controlling process.

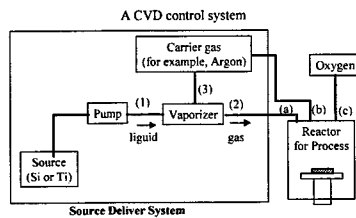


Dept. of Electrical and Computer Engineering

Application of Gaze Detection Technology to Process Controls

Application of Gaze Detection Technology to Process Controls - (3)

1 Examples



The pipes (a, b, c) carrying oxygen, argon and vaporizer need to be opened at the same time, or if the pump is not working well, then the pipes (1), (2) and (3) need to be closed at the same time.
→ The pipes(a, b) or (1, 2) are manipulated by both hands and the pipe(c) or (3) is handled by gaze detection technique

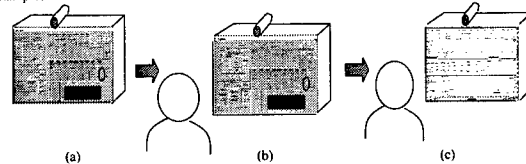
Dept. of Electrical and Computer Engineering

Application of Gaze Detection Technology to Process Controls

Application of Gaze Detection Technology to Process Controls - (4)

1 In its second function, gaze detection technology can store the process information during the user's absence and replay it when he returns

1 Examples



- (a) Store the process information(pressures and temperatures in the reactor) which happened during the user's absence
- (b) Recognize whether he returned or not
- (c) Replay the stored information when he returned

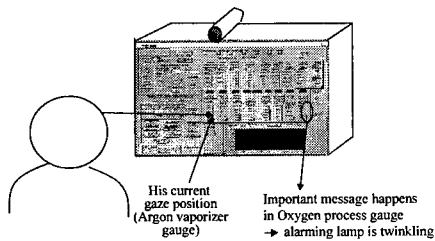
Dept. of Electrical and Computer Engineering

Application of Gaze Detection Technology to Process Controls

Application of Gaze Detection Technology to Process Controls - (5)

1 In its third function, gaze detection technology can lead a user's attraction to important spots

1 Examples

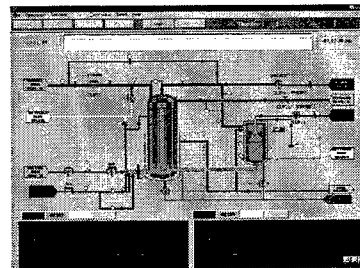


Dept. of Electrical and Computer Engineering

Application of Gaze Detection Technology to Process Controls

Application of Gaze Detection Technology to Process Controls - (6)

1 A user interface of ready-made Syn-gas Composition program made by Chemical Engineering of Yonsei Univ.



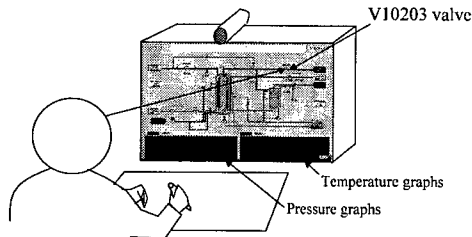
Dept. of Electrical and Computer Engineering

Application of Gaze Detection Technology to Process Controls

Application of Gaze Detection Technology to Process Controls - (7)

In its first function, gaze detection technology enables a user to control another process (V10203 valve opening), even when both hands are busy in writing a note about the process information (temperature and pressure).

Examples



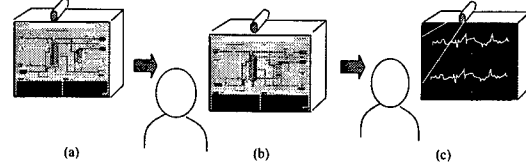
Dept. of Electrical and Computer Engineering

Application of Gaze Detection Technology to Process Controls

Application of Gaze Detection Technology to Process Controls - (8)

In its second function, gaze detection technology can store the process information during the user's absence and replay it when he returns.

Examples



- (a) Store the process information (pressures and temperatures in the Syn-gas Composition Reactor) which happened during the user's absence
- (b) Recognize whether he returned or not
- (c) Replay the stored information when he returned

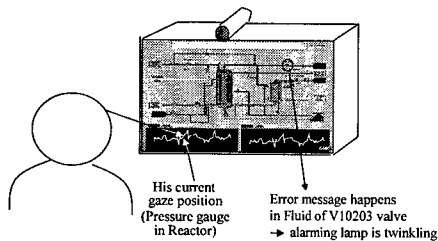
Dept. of Electrical and Computer Engineering

Application of Gaze Detection Technology to Process Controls

Application of Gaze Detection Technology to Process Controls - (9)

In its third function, gaze detection technology can lead a user's attraction to important spots.

Examples



Dept. of Electrical and Computer Engineering

Application of Gaze Detection Technology to Process Controls

Other Applications - (1)

Interface system for the handicapped



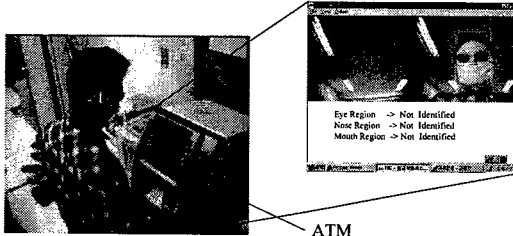
Dept. of Electrical and Computer Engineering

Other Applications

Other Applications - (2)

Face Identification System in ATM (Automatic Teller Machine)

When a dubious user uses ATM, who can not be identified by Face Identification System, he is rejected for transaction until his face is identified.



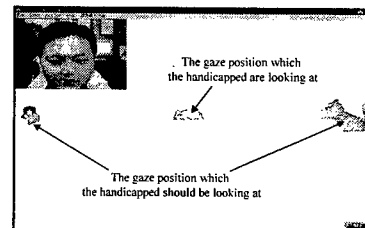
Dept. of Electrical and Computer Engineering

Other Applications

Other Applications - (3)

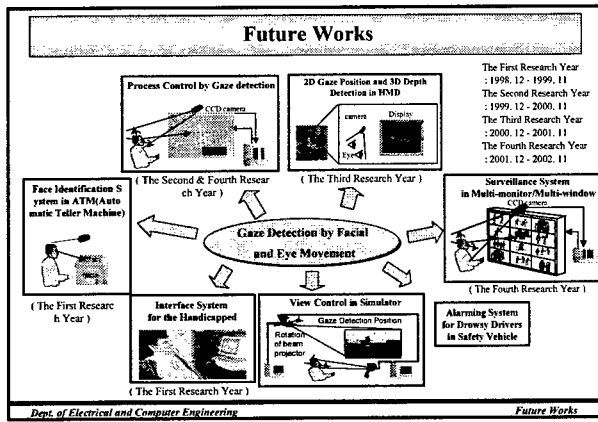
The rehabilitation therapy system for the handicapped

→ This system measures the gazing time and accuracy of the handicapped by gaze detection technique and uses it for rehabilitation therapy.



Dept. of Electrical and Computer Engineering

Other Applications



The Third Eye Cameras - Dynamic and Static Hyperspectrum Imaging

Yoshiyasu Takefuji

Faculty of Environmental Information, Keio University,
5322 Endo, Fujisawa, 2520816 Japan
Email: takefuji@sfc.keio.ac.jp

The goal of this paper is to introduce our third-eye camera projects using dynamic and static hyperspectrum imaging with integration of several cameras including a new super harp camera, IR camera. The third-eye camera is able to capture value-added images of objects with a wide range of spectrum from radio wave, long IR, IR, visible light, ultraviolet, or X-ray by neural feature extraction. Static hyperspectrum imaging can be used for inspection of buildings and bridges, investigating the status of tiles, underground and moss, investigation of coral, measuring gases for health of plants, measuring the taste of food, discovering counterfeit money, minute blood vessel and so on. Artificial retina camera is also introduced as dynamic spectrum imaging for many applications including optical flow imaging, edging, and so on. A real-time gesture interpreter using the retina camera will be addressed.

The Third-EYE Project hyperspectral computing

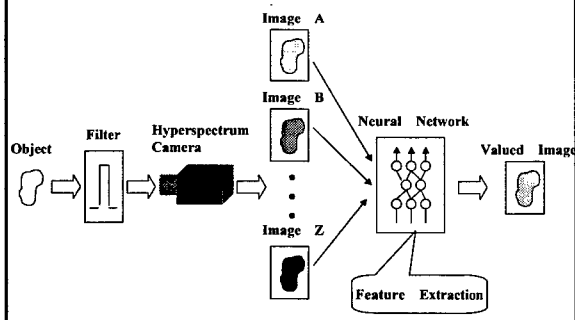
Keio University

Yoshiyasu Takefuji

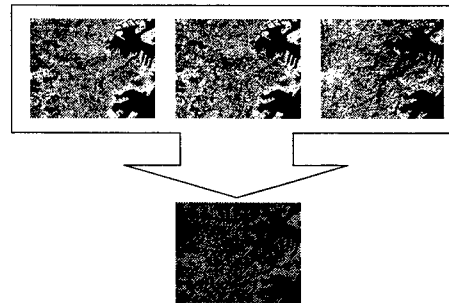
The goal of hyperspectral imaging:

- Valuable features can be captured from the wide range of hyperspectral data including radio wave, long infrared, infrared, visible light, ultraviolet, X-ray,...

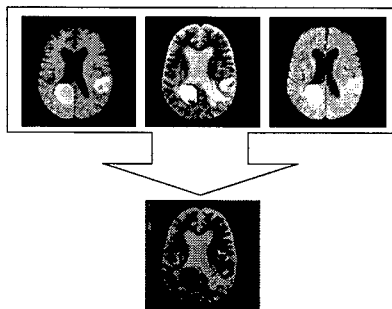
Overview of The Third Eye system



Multispectral Remote-Sensing Imaging



Multispectral MRI Imaging



The role of UltraViolet



Human's Eye



Bee Cam

©NHK

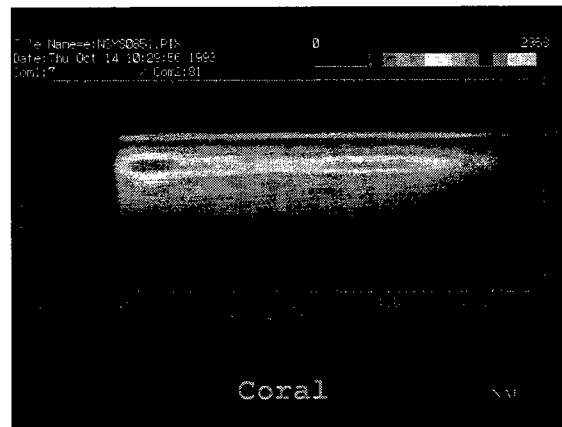
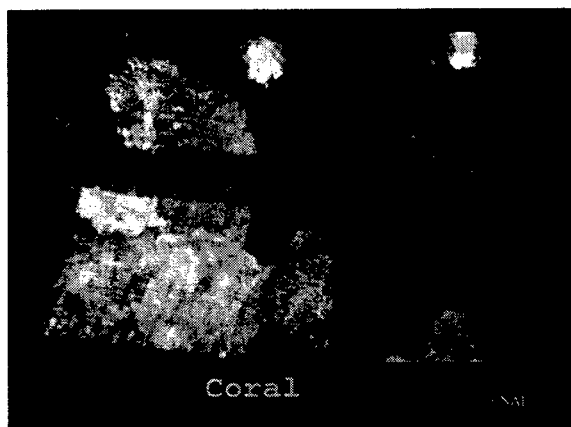
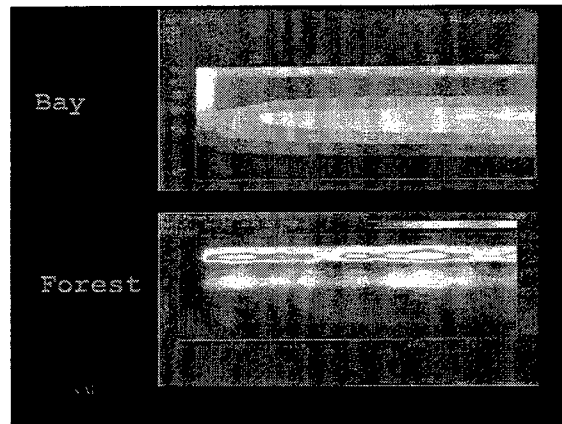
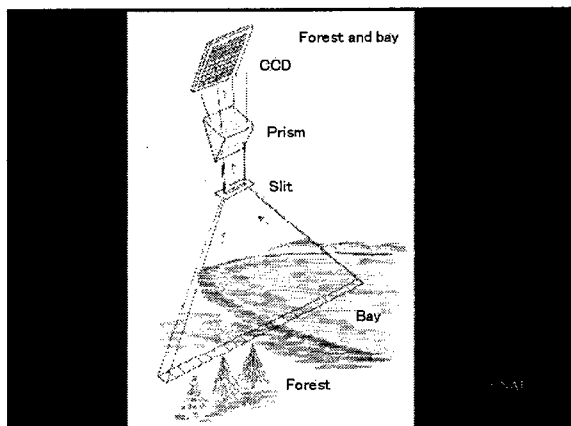
NASA space shuttle tile inspection



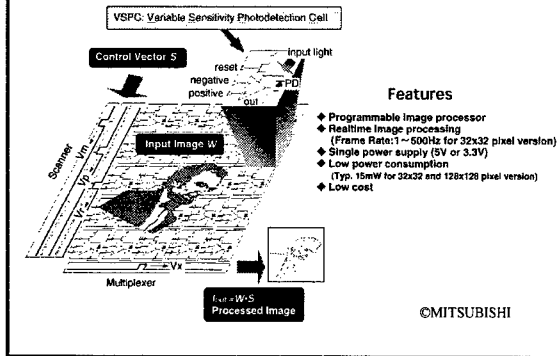
©NASA

Hyperspectral applications

- tile or wall
- moss or river
- gases from plants
- coral
- printing
- health of earth
- taste
- art paints
- rusts



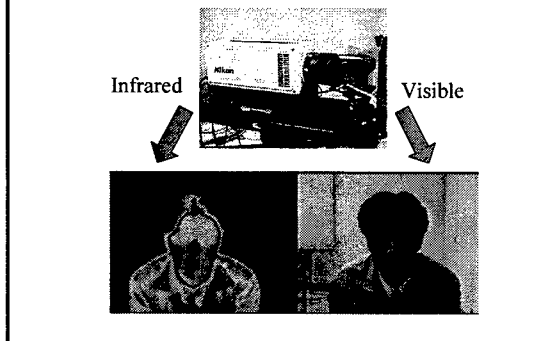
Configuration of Mitsubishi CMOS Image Sensor



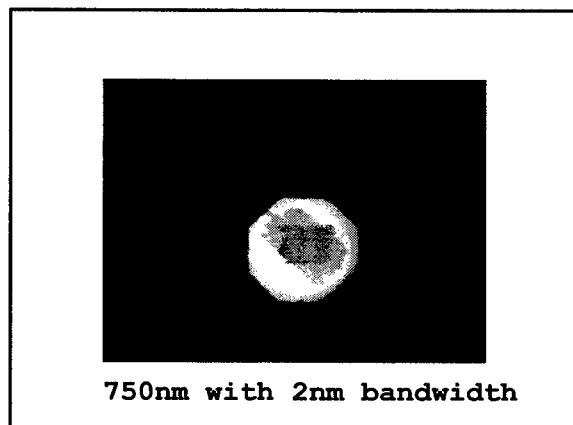
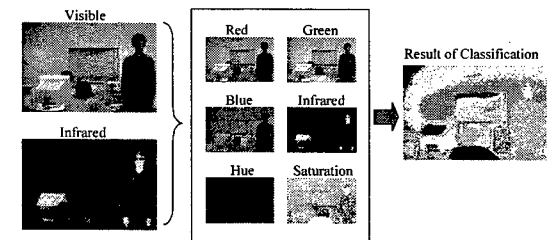
Algorithm Lineup for Mitsubishi CMOS Image Sensor

Algorithm	Image Movements Recognition Algorithm	Image Shapes Recognition Algorithm	Character Recognition Algorithm	Position Recognition Algorithm
Examples				
Principle	Optical Flow	Detection of Center of Gravity	Pattern Matching	Image Projection
Applications	<ul style="list-style-type: none"> • Body Action Recognition • Gesture Recognition • Person Detection 	<ul style="list-style-type: none"> • Gesture Recognition • Person and Car Recognition 	<ul style="list-style-type: none"> • Numerals Recognition • Bar Code Recognition 	<ul style="list-style-type: none"> • Coordinates Detection

IR+CCD Cameras

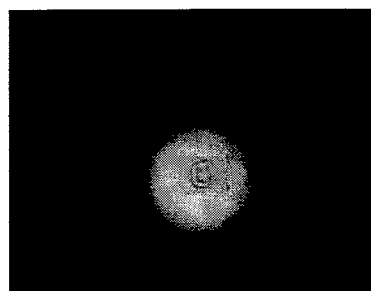


Human's Feature Extraction



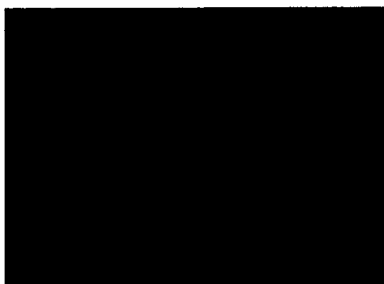


800nm with 2nm bandwidth



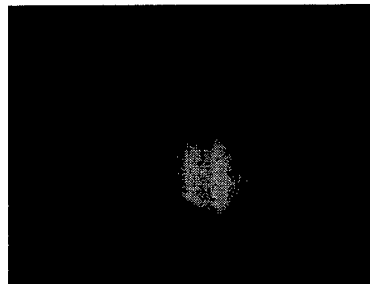
850nm with 2nm bandwidth

670nm



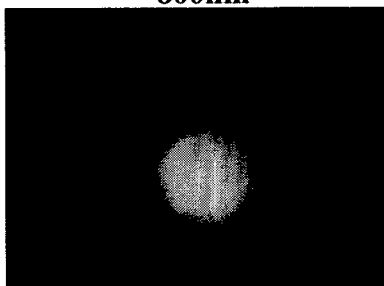
Organic/non-organic onion

700nm

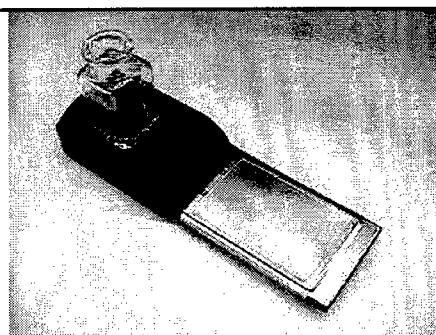


Organic/non-organic onion

800nm



Organic/non-organic onion



Field-Spectrometer

©Nomadics

The Third Eye Approach to Innovative Designs and Applications: Human Recognition System by Nonlinear Oscillations

Souichi Oka, Yoshiyasu Takefuji, William Huang

*Graduate School of Media and Governance, Keio University
5322 Endo, Fujisawa, Japan 252-0816,

Email: takefuji@sfc.keio.ac.jp, souichi@sfc.keio.ac.jp, william@sfc.keio.ac.jp

The field of remote sensing and sensor technology has undergone tremendous development in the past decades. Sensor technologies of all kinds such as electro-optics, acoustic, active/passive UV to LWIR, ground-penetrating radar, passive mm wavelength, x-ray tomography, neutron activation imaging, multi-spectral, hyper-spectral, and ultra-spectral imaging, provide valuable images that normal CCD cameras cannot offer. By combining algorithms and images taken by sensors at different ranges of the electromagnetic spectrum, we will be able to extract valuable images automatically. By using multi-spectral images and processing them with neural network computing, our "Third Eye" team is able to extract human face features from those images. In this paper, we will present an application for detecting human facial parts, images taken by different imaging systems and sensors, and the current status of image processing applications.

Light is a form of electromagnetic radiation. Other forms of electromagnetic radiation include radio waves, microwaves, infrared radiation, ultraviolet rays, X-rays, and gamma rays. All of these are known collectively as the electromagnetic spectrum. The rainbow of colors that we see in visible light represents only a very small portion of the electromagnetic spectrum. The electromagnetic spectrum ranges from gamma rays to radio waves with everything else in between. The reflectance spectra of most materials on the Earth's surface contain characteristics or diagnostic absorption features. Remote sensors, such as hyper-spectral imager developed by TRW, capable of acquiring complete reflectance spectra over large areas offer a powerful tool for study of the Earth and the environment. With a CCD camera, we are only able to see the visible part of the electromagnetic spectrum that ranges from 0.4 μm to 0.7 μm of the spectrum, thus, any range below or beyond the visible part of the EMS are invisible to human eyes. However, by combining images taken at different part of the spectrum and processing them, we will be able to reveal many valuable images. Unique algorithms used in conjunction with different imaging systems and sensors will allow exploration of a large variety of basic problems in fields such as earth science, vegetation studies, geology, semiconductor manufacturing, biology, biological imaging, material processing, environment, medical imaging, ground target detection, chemical identification, and inspection.

In the field of machine vision and visual inspection, a company called SRI International has pioneered, designed, and developed an automatic visual inspection system that has applied widely in inspection of defects, cracks, or structural failures in different area such as on NASA shuttle tile inspection. There are over 20,000 thermal tiles that cover the outer surfaces of the shuttle. Each of these tiles is life-critical because a failure can result in the shuttle's aluminum skin overheating on reentry, thus, causing a catastrophic system failure. Consequently, each tile is manually inspected after each flight to ensure that there are no cracks, dents, or other structural failures that require rework. To reduce human error and inspection time, SRI International has developed an imaging system and inspection algorithm to detect damages to the thermal tiles. The picture on the left is the preflight image, and the picture on the right is the postflight image. The right image shows a simulated defect identified by their algorithms. This system demonstrates the fact that by building or manipulating different sensors and developing specific algorithms, we will be able to produce incredible images which will facilitate us in defect inspection or manufacturing that will ultimately save time, money, and human lives.

**THE THIRD EYE APPROACH TO
INNOVATIVE DESIGNS AND APPLICATIONS :
HUMAN RECOGNITION SYSTEM BY
NONLINEAR OSCILLATIONS**



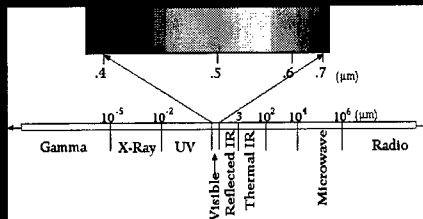
What can you see ?

Keio University
Souichi Oka

Contents

- Introduction
- Algorithm
- Simulation
- Conclusion

Applications using Hyper-Spectral Camera



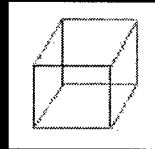
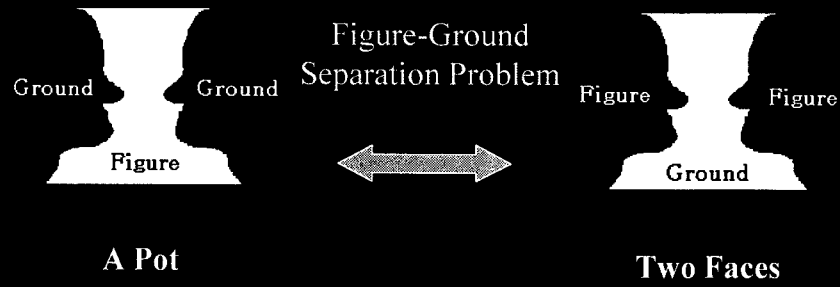
Human Recognition Problem



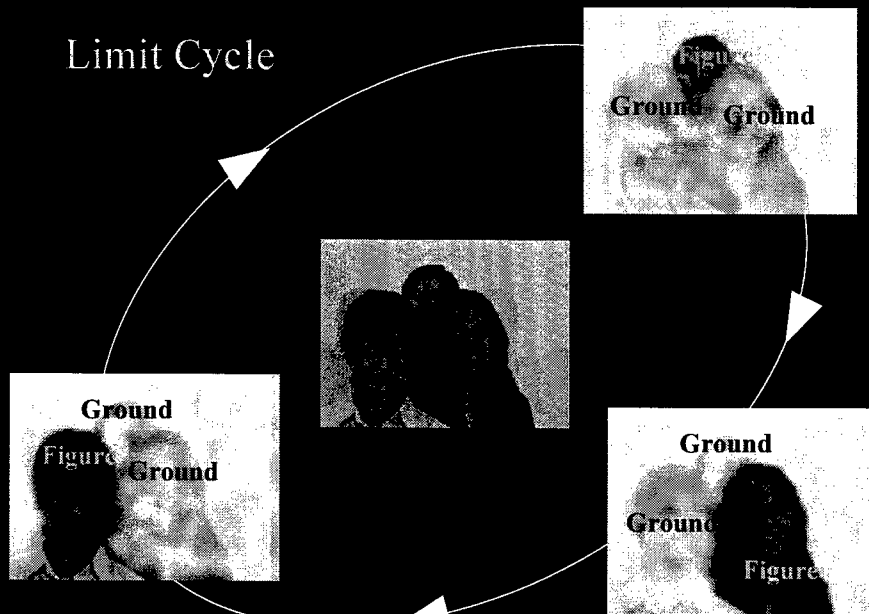
The goal of this study is to discriminate several persons by the proposed neural network.

This system groups facial parts to detect a face.

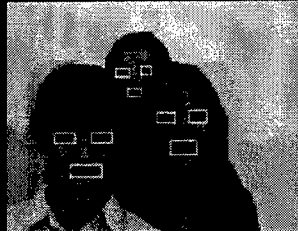
Synchronized Brain Hypothesis



Limit Cycle



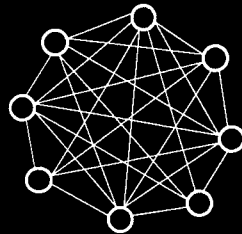
Grouping Facial Parts Model



Extract facial parts including eye or mouth.



Extracted facial parts are fed to recurrent neural network.



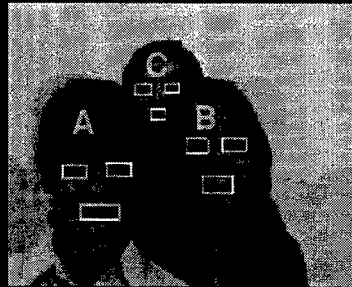
Our idea:

- 1 : A principle of similar figure
- 2 : Nonlinear Oscillatory neurons

Motion Equations

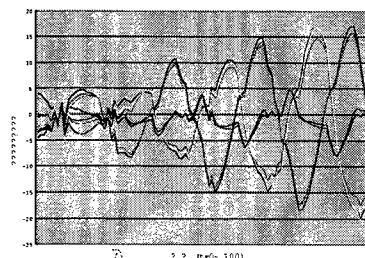
$$\begin{aligned}
 dX_i = & -A(\sum V_i \cdot eye(i) - 2) \quad \text{Eyes} \\
 & -A(\sum V_i \cdot mouse(i) - 1) \quad \text{Mouth} \\
 & -C \sum \sum \sum |real_angle(i, j, k) - ideal_angle(i, j, k)| \quad \text{Angle} \\
 & -D \sum \sum \sum \sum V_i V_j V_m V_n \left| \frac{real_length(i, j)}{ideal_length(i, j)} - \frac{real_length(m, n)}{ideal_length(m, n)} \right| \\
 & + EV_p V_q \frac{1}{2\pi} e^{-\frac{1}{2}(((ideal_x - x(i))^2 + (ideal_y - y(i))^2)/\sigma^2)} \quad \text{Ratio of scale} \\
 & -f(Y_i) \quad \text{Feedback from the inhibitory term} \\
 & + \int_0^t X_i e^{\tau-t} d\tau \quad \text{Fatigue of the neuron} \\
 dY_i = & -\frac{Y_i}{\tau_y} + f_y \left(-\frac{T_w Y_i}{\bar{Y}} + \frac{T_x X_i}{\bar{X}} \right) \quad \text{Inhibitory term}
 \end{aligned}$$

Sample Target Image

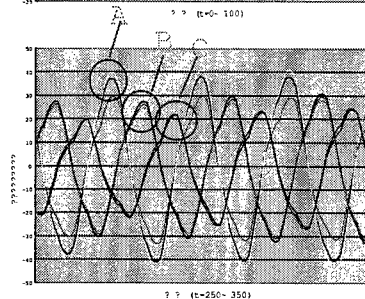


Nine facial parts are fed to the proposed neural network.

From Chaos to Order



At first, the dynamics is instability, and synchronized gradually. ($t=0 \sim 100$)



Once the dynamics converges to the limit cycle, three persons are detected one after another. ($t=200 \sim 300$).

Conclusion

- Our system can flexibly discriminate several human faces in a time-varying image.
- The normalization problem is resolved about scale, rotation, and shift of an object.

Intelligent Rate Control for MPEG-4 Coders

Gwang Hoon Park, Jae Hyung Park and Yoon Jin Lee

Department of Computer Science, Yonsei University,
234 Maji, Heungup, Wonju, Kwangwon, Korea
Email: ghpark@dragon.yonsei.ac.kr

The multimedia technology, which will lead in the 21st century, will depend on how to manipulate visual information efficiently, and will be focused on the interactive visual information exchange and user interactions. Among one of those technologies, MPEG-4 codec can multiplex a set of independently-coded, arbitrarily-shaped video objects and transmit through either fixed or variable rate channels such as internet, wireless or satellite communications. Some quality control algorithms should support the encoding of visual objects to obtain and maintain best picture quality under the constraints of the quality requirements and channel environments. Especially, MPEG-4 Codec should be robust with respect to rapid changes in size and shape of the objects.

This paper focuses on the design of the intelligent rate control algorithms via introducing quadratic neural networks and evaluating data-driven pattern analysis rather than rate-distortion mathematical models. According to data-driven pattern analysis, it is found that several new variables such as motion vectors are required to control near-optimally for transmitting best quality of moving pictures in real-time. And we also found that the simplified mathematical rate distortion models, which are now widely used, could not support the control mechanism enough. Therefore, quadratic neural-net using density estimation and randomizing pattern space, called Density-based Random-Vector Functional-link Net is introduced to control the picture quality optimally. The proposed algorithm is tuned to obtain near-optimal picture quality of QCIF (176X144 pixels) format video streams for low bitrate transmission, storage/retrieval. The experimental work will be presented based on the recent MPEG-4 video coder specification, apart from making intelligent rate control algorithm. The comparisons between the results of intelligent control and conventional rate control will be presented in this paper.

Intelligent Rate Control for MPEG-4 Coders⁺

Gwang Hoon Park*, Jae Hyung Park*,
Yoon Jin Lee* and Steven R. LeClair**

* Department of Computer Science, Yonsei University, Wonju, Korea

** Materials Process Design Branch, Manufacturing Technology Division, Air Force Research Laboratory, U.S.A.

⁺ This work was supported in part by U.S. Air Force Office of Scientific Research (AOARD-98-4011)

G. H. Park, J. H. Park, Y. J. Lee, and S. R. LeClair, Yonsei University



Summary

- MPEG-4 codec can multiplex a set of independently-coded, arbitrarily-shaped video objects and transmit it through either fixed or variable rate channels such as internet, wireless or satellite communications.
- Quality control algorithms should support the encoding process of the visual objects to obtain and maintain best picture quality under the constraints of the quality requirements and channel environments.
- This paper focuses on the design of the intelligent rate control algorithm via introducing global rate distortion (RD) model constructed by quadratic neural network, by evaluating data-driven pattern analysis rather than rate-distortion mathematical models.
- Proposed global RD model is very useful in case the characteristics of the video sequences are rapidly varying.
- The regression based mathematical model may not support rapidly changing environments, because it requires the time to stabilize to generate appropriate Q steps.
- Intelligent rate control based on the global RD model can generate appropriate Q steps immediately in feedforward manner.
- The performances of the proposed algorithm are superior than those of the MPEG-4 VM5+ rate control based on the regression process, in comparison with the average bits per frame to satisfy the channel constraints, encoded peak Signal-to-Noise Ratio (PSNR), and the number of frame skips.

G. H. Park, J. H. Park, Y. J. Lee, and S. R. LeClair, Yonsei University



Introduction

- **Image Data Compression**
 - Reduce Information of the Image that is not perceptible by Humans.
 - MPEG uses a hybrid coding methodology that is divided into two parts:
 - the algorithm for reduction of temporal redundancy between adjacent frames
 - by motion estimation and compensation technique (ME/MC)
 - the algorithm for reduction of spatial redundancy .
 - using DCT (Discrete Cosine Transform) within image frame
- The quality of MPEG video bitstream can be improved by the various encoding algorithms flexibly designed by the experts.
 - Pre/Post-processing, Motion Estimation, Rate Control algorithms
(can be flexibly designed for specific coding environment.)
- To transmit the huge amount of video information through a bandwidth constrained channel such as mobile picture communications or internet TV broadcasting on the ISDN or PSTN networks
 - Some quality control algorithms should support the encoding process of the visual objects to obtain and maintain best picture quality under the constraints of the quality requirements and channel environments.
- We will focus on the MPEG-4 video compression algorithm and
Propose a new intelligent technique for MPEG-4 rate control.

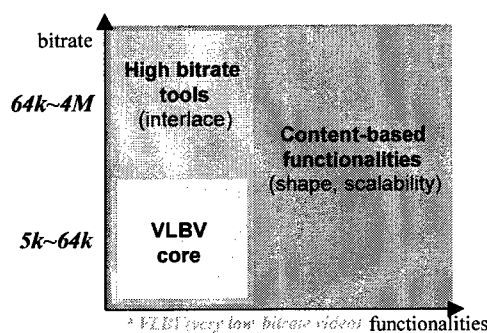
© 1998, Beijing University of Aeronautics and Astronautics, Tsinghua University



MPEG-4 Visual

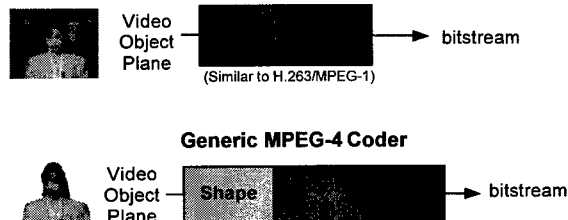
Natural Video

VOP based Coding
Content-based coding
Content-based scalability
Error resilience
Efficient compression
Efficient random access
Extended manipulation



Synthetic Video

Facial & body animation
Object manipulation



© 1998, Beijing University of Aeronautics and Astronautics, Tsinghua University



Related Theory

According to information theory, two problems are stated:

source coding (what information should be sent)

channel coding problem (how should it be sent).

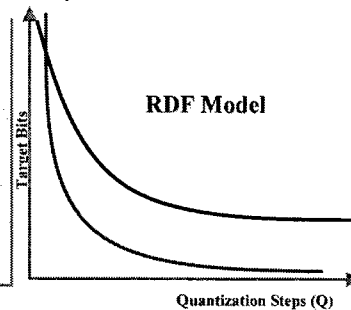
Rate Distortion Theory (RDT) is directly related to the source coding problem
(lossy image data compression)

The RDF model has been considered as a good choice to represent relations between quantizing distortions and encoder output rates and thus has been used in wide range.

The key factor in RDT is the rate distortion function (RDF) $R(D)$, which represents the lower bound on the rate:

If a certain channel capacity C is given, RDF is used to find

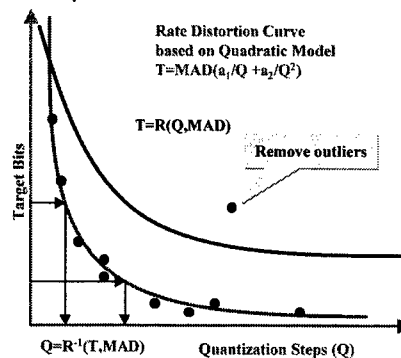
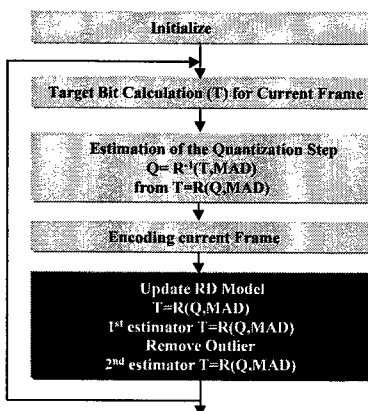
the necessary minimum average distortion D_{ave} so that the condition for error-free transmission $R(D_{ave}) < C$ is achieved.



© 2004 Intel Corporation. All rights reserved. Intel Corporation. Confidential.

INTEL

MPEG-4 VM Rate Control



Choice of quantizer steps at the encoder plays a key role in determining the actually encoded bitrate and the quality of the transmitted video scenes.

The recommended rate control algorithm in MPEG (3 steps)

Bit allocation : Past bit usage and quantizer steps are used to estimate the relative complexity of the picture and thereby determine the target bit rate for the present picture

Rate control : A reference quantizer step is determined by evaluating a virtual buffer status and the difference between the target bit rate and the rate that is already consumed till now

Adaptive quantization based on the mathematical model : Regression based on mathematical model is carried out to decide actual quantizer for the present frame or macroblocks.

© 2004 Intel Corporation. All rights reserved. Intel Corporation. Confidential.

INTEL

MPEG-4 VM Rate Control Algorithm

- Quadratic rate distortion model is used to estimate the rate distortion curve to evaluate the target bit rate before performing the actual encoding
- Recommended MPEG-4 VM rate control algorithm

Initialization

Computation of the target bit rate before encoding.

- ✧ The target bit rate is computed based on the bits available and the last encoded frame bits. If the last frame is complex and uses excessive bits, more bits should be assigned to this frame. However, there are fewer bits left for encoding. Thus, fewer bits can be assigned to this frame.
- ✧ A lower bound of target bit rate ($F/30$) is used so that the minimal quality is guaranteed (F : total target bits per second).
- ✧ The target bit rate is adjusted according to the buffer status to prevent both overflow and underflow.

Computation of the quantization parameter (Q) before encoding

- ✧ Q is solved based on the model parameters, a_1 and a_2 .
- ✧ Q is clipped between 1 and 31.
- ✧ Q is limited to vary within 25% of the previous Q to maintain a variable bit rate (VBR) quality.

Encoding current frame

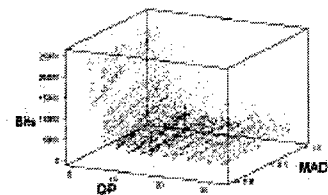
Model parameters are updated based on the encoding results of the current frame.

- ✧ The rate distortion model is updated based on the encoding results of the current frame. The bits used for the header and the motion vectors are deducted since they are not related to Q.
- ✧ The data points are selected using a window whose size depends on the change in complexity. If the complexity changes significantly, a smaller window with more recent data points is used.
- ✧ The model is calibrated again by rejecting the outlier data points. The rejection criterion is the data point is discarded when the prediction error is more than one standard deviation.
- ✧ The next frame is skipped if the current buffer status is above 80 %.

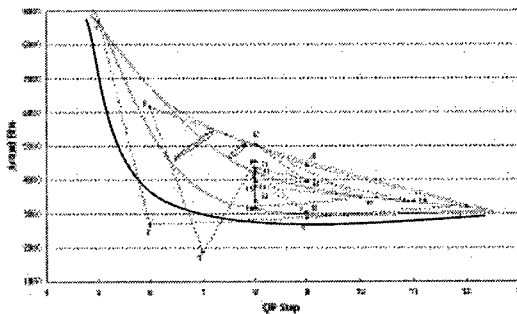
G. H. Park, Image Processing with Computer Science in Video Technology



Actually Generated Rate Distortion Curves



(a) Relation between actually encoded bits, MAD and Q taken from several QCIF image sequences



(b) Example of RD curves based on actually generated relations between encoded bits and Q steps for twenty image frames.

- Possible Global RD model.
- However we may not construct the global RD model by using only T, MAD and Q, because there are too many variations to formulate the correct global RD model.

For Q (QP step in the figure)= 6,
 about 2500 bits are generated at encoding of the 2nd frame
 about 6500 bits need to be coded at 6th image frame.
 Frame skip may be occurred while encoding 6th frame
 (about 4000 bits are excessively generated, therefore buffer may be full)
 by using RD curve predicted by regression of the group of the points 1,2,3, and 4 (point 5 may be eliminated by removal process in constructing mathematical RD model)

G. H. Park, Image Processing with Computer Science in Video Technology



Global Rate Distortion Model

- Need to generate fast-moving appropriate RD curves to adapt rapidly varying image characteristics to obtain nearest encoded bits to the target bits assigned.
- To formulate global RD model, several additional parameters are needed in addition to T, MAD and Q.
- Not to have additional processes to find parameters of the global RD model, find parameters in the necessary processes of the encoding procedure.
- Before rate coding is performed, motion estimation and compensation (ME/MC) are carried out to reduce temporal redundancies.
- If the motion varied too much or objects in image frame is just disappeared or created, just intra block coding is performed in the spatial domain.
- If the intra block coding is performed in the macroblock, large bits are required to encode corresponding macroblock.
- For the motion part, we divide it into three parameters, such as MV_1 , MV_{2-9} , MV_{10+} (in case motion vector range is 16 in QCIF format image sequences).
- To simplify the algorithm, only histogram of the moved macroblocks in specific motion range is used. MV_1 is the number of macroblocks that moved within 1 pixel range, MV_{2-9} is the number of macroblocks that moved from 2 to 9 pixel range
- MV_1 , MV_{2-9} , MV_{10+} , IntraMB and MAD are already known before rate control is carried out in the encoding process.
- It is very difficult to formulate global RD model based on mathematical analysis, therefore we use the neural network to construct it, based on data-driven pattern analysis.

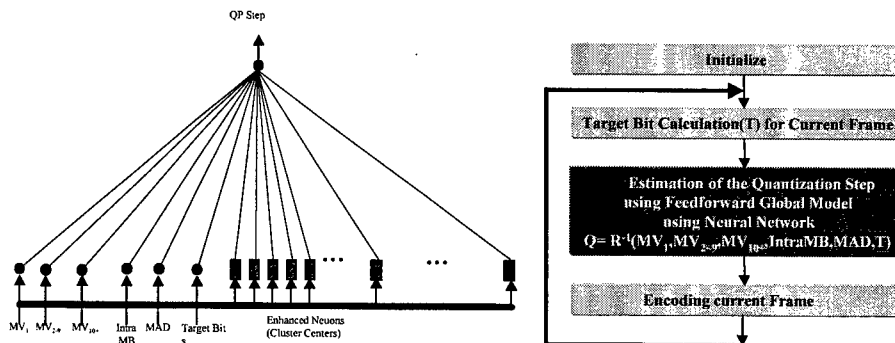
$$T = R(MV_1, MV_{2-9}, MV_{10+}, IntraMB, MAD, Q)$$

$$Q = R^{-1}(MV_1, MV_{2-9}, MV_{10+}, IntraMB, MAD, T)$$

Y. G. Hu, P. Hui, Institute of Information Science, Tsinghua University

CMREC

Feedforward Intelligent Rate Control



Radial Basis Function Neural Network is used to formulate

Global Rate Distortion Model (Inverse Problem)

Quadratic RBFNN is used

Conjugate Gradient (CG) method is used for Training (N+6 iterations)

N Cluster Centers are obtained using K-means Algorithm

(Classifications of the variations of Input Space ($MV_1, MV_{2-9}, MV_{10+}, IntraMB, MAD, T$))

Y. G. Hu, P. Hui, Institute of Information Science, Tsinghua University

CMREC

Intelligent Rate Control Algorithm

Initialization

- Load the values of weights and centers of the radial basis function neural network.

Computation of the target bit rate before encoding (same process used in MPEG-4 VM).

- ✧ The target bit rate is computed based on the bits available and the last encoded frame bits. If the last frame is complex and uses excessive bits, more bits should be assigned to this frame. However, there are fewer bits left for encoding. Thus, fewer bits can be assigned to this frame. A weighted average reflects a compromise of these two factors.
- ✧ A lower bound of target bit rate ($F/30$) is used so that the minimal quality is guaranteed (F : total target bits per second).
- ✧ The target bit rate is adjusted according to the buffer status to prevent both overflow and underflow.

Computation of the quantization parameter (Q) via feedforward manner.

- ✧ Estimation of the quantization step using global RD model generated by RBFNN.
- ✧ Q is clipped between 1 and 31.

Encoding current image frame.

- ✧ The next frame is skipped if the current buffer status is above 80 %.

© 2001, Pattern Image Processing Lab, Computer Science, Tsinghua University



Simulation Results for MPEG-4 Reference Images



News



Akiyo



Hall



Silent



Container

MPEG-4 reference image sequences

Comparison of the performances generated by the rate controls of the MPEG-4 VM and the intelligent rate control for the MPEG-4 reference image sequences (QCIF, 48kbps, 10frames/sec). No frame skip is occurred.

48kbps, 10 frames/sec		Training			Test	
Sequences		News	Akiyo	Hall	Silent	Container
MPEG-4 VM	Ave. bits / VOP	4799.43	4796.95	4816.28	4857.67	4796.25
	PSNR (Y)	33.39	40.37	36.68	33.71	34.93
	PSNR (U)	37.41	42.38	39.30	37.34	40.19
	PSNR (V)	38.29	43.71	41.25	38.56	39.69
Intelligent Rate Control	Ave. bits / VOP	4799.26	4821.48	4802.45	4822.27	4796.50
	PSNR (Y)	33.46	40.39	36.79	33.85	35.02
	PSNR (U)	37.39	42.49	39.33	37.53	40.22
	PSNR (V)	38.21	43.84	41.28	38.70	39.76

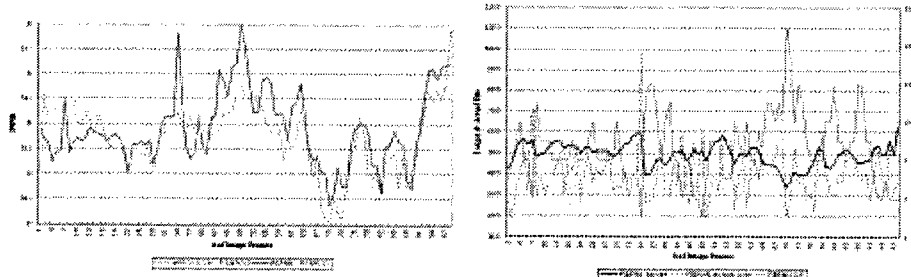
Performances of the intelligent rate control are around 0.1 dB better in PSNR's than those of the MPEG-4 VM rate control.

Even the intelligent rate control produce less bits than the MPEG-4 VM, in average bits per frame.

© 2001, Pattern Image Processing Lab, Computer Science, Tsinghua University



Comparisons of the PSNR's encoded by MPEG-4 VM rate control and Intelligent Rate Control



(a) Comparison of the PSNR's (in dB) of the luminance parts (Y) encoded by model based MPEG-4 VM rate control and the RBFNN based intelligent rate control ('Silent' Sequence).

(b) Comparisons between target bits and actually encoded bits by RBFNN based intelligent rate control and corresponding Q steps.

Intelligent rate control quickly stabilizes Q steps concurrently with connected to the variations of the image characteristics.

© 1998, Intel Corporation. All rights reserved. Intel Corporation. Intel Corporation.

INTEL CORPORATION

Simulation Results for Pulsed-Laser-Deposition Plume Video Sequences



Pulsed-Laser-Deposition plume video sequences

Comparisons of the simulation results encoded by the rate controls of the MPEG-4 VM and the intelligent rate control for the Pulsed-Laser-Deposition plume video sequence (QCIF Format, 30 frames/second, 112kbps and 256kbps)

30 frames/sec		112kbps		256kbps	
500 image frames		Training	Test	Training	Test
MPEG-4 VM Rate Control	Ave. bits / VOP	5450.18	5529.00	9906.32	10025.74
	PSNR (Y)	34.55	34.65	37.39	37.43
	PSNR (U)	38.92	38.97	40.40	40.45
	PSNR (V)	39.98	40.06	42.05	42.11
	Frame skip	155	160	66	71
Intelligent Rate Control	Ave. bits / VOP	4403.33	4393.72	8515.94	8517.61
	PSNR (Y)	33.76	33.72	36.64	36.65
	PSNR (U)	38.65	38.56	39.95	39.98
	PSNR (V)	39.41	39.37	41.32	41.36
	Frame skip	73	72	0	0

At 112kbps:

32% (160/500) of the frames are skipped in the MPEG-4 rate control. 14.4% (72/500) frames are skipped in the intelligent rate control.

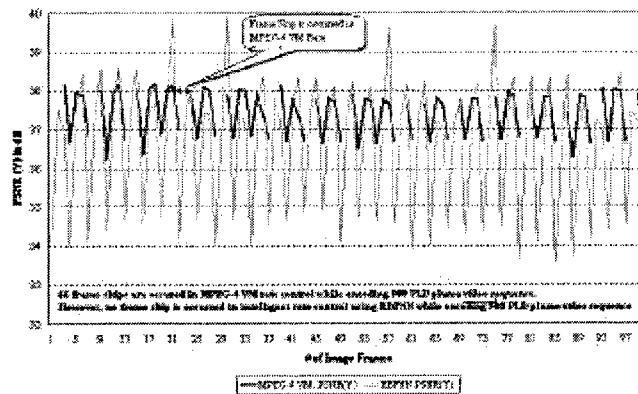
At 256kbps:

71 frames are skipped in the MPEG-4 rate control. Every frame is properly coded by the intelligent rate control algorithm.

© 1998, Intel Corporation. All rights reserved. Intel Corporation. Intel Corporation.

INTEL CORPORATION

Simulation Results for Pulsed-Laser-Deposition Plume Video Sequences



Intelligent rate control encodes every plume appeared images.

MPEG-4 VM can not encode those image scenes due to less accuracy of the regression based mathematical model, therefore followed by the frame skips.

Comparison of the PSNR's (in dB) of the luminance parts encoded by mathematical model based MPEG-4 VM rate control and the RBFNN based intelligent rate control for the Pulsed-Laser-Deposition plume video sequences: 71 frame skips are occurred in MPEG-4 VM rate control, however no frame skip is occurred in the intelligent rate control while encoding 30 frames/second 500 PLD plume video sequence at 256kbps.

C. G. H. Park, Image Processing Lab, Chungnam National University, Taejeon, Korea

CONFERENCE

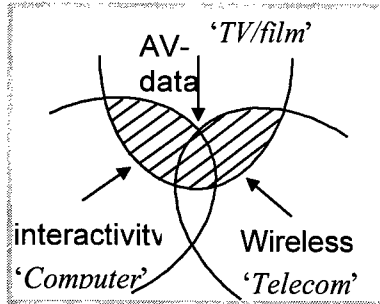
Conclusions

- Introduce feedforward global Rate Distortion model using Radial Basis Function neural Network.
- The performances of the proposed algorithm are superior than those of the MPEG-4 VM rate control algorithm based on mathematical RD model
 - in comparisons with the average bits per frame to satisfy the channel constraints, encoded PSNR's, number of frame skips.
- Proposed global RD model is very useful in case the characteristics of the video sequences are rapidly varying.
 - The regression based mathematical model may not support rapidly changing environments, because it requires the time to stabilize to generate appropriate Q steps.
 - Intelligent rate control based on the global RD model can generate appropriate Q steps immediately in feedforward manner, because neural network already trained many cases of the variations of the image characteristics.
- The proposed intelligent rate control can be usefully used in the case of
 - the generic video phones,
 - broadcasting of the sports games via internet,
 - special purpose video scenes like PLD plume video sequences
 because those video sequences have similar characteristics, therefore the possible variations of the image characteristics can be trained before the encoding current video scenes.

C. G. H. Park, Image Processing Lab, Chungnam National University, Taejeon, Korea

CONFERENCE

Multimedia Trends & MPEG-4 Focus



Three Major Trend (Multimedia)

- trend towards **wireless communication**
- trend towards **interactive Computer Applications**
- trend towards **integration of A/V data into an ever increasing # of applications (TV/Film)**

- MPEG-4 : Applications in the Shaded Area
- Focus of MPEG-4
 - Content-based Interactivity
 - High Compression (Low Bitrate)
 - Universal Accessibility over wide range of storage and transmission media
 - High degree of Flexibility and Extensibility (MSDL)

© 1998, Intel, Image Processing Lab, Computer Science, Texas A&M University

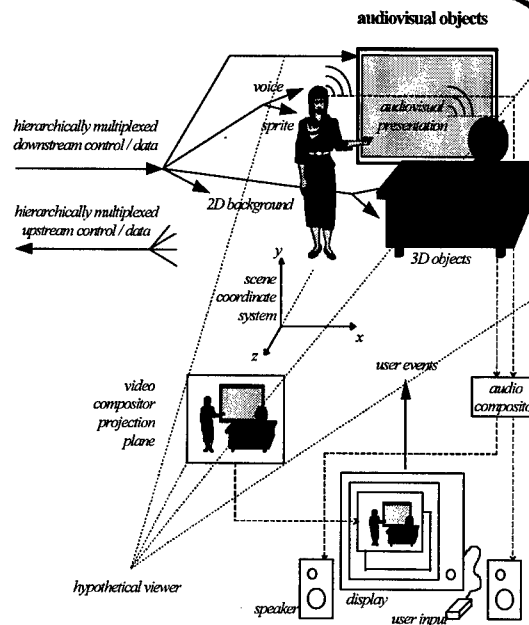
MPEG-4

MPEG-4 Model

AVO: audio/visual Object
 2-d background (sprite)
 2-d/ 3-d objects
 talking person
 voice of the person
 talking head
 animated human body
 Text stream (TTS)

Compound AVO objects
 ex: visual + speech, audio
 by the compositor

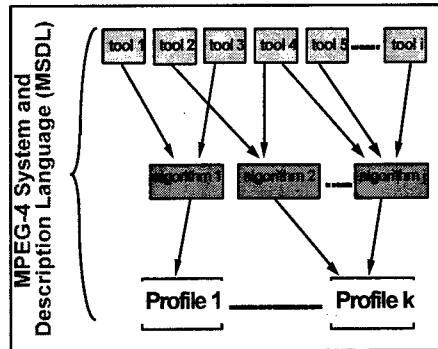
User Interaction
 manipulate meaningful
 objects



© 1998, Intel, Image Processing Lab, Computer Science, Texas A&M University

MPEG-4

MPEG-4 Structure



- Tools

- Shape Coding
- Motion Est./Comp.
- Texture Coding
- etc.

- Algorithms

Collection of the tools

- MPEG-1 Video, MPEG-2 System, Video, Audio, etc.

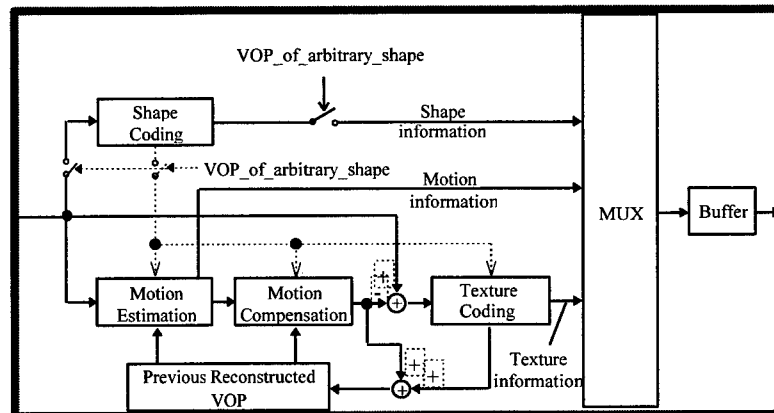
- Profiles

- MPEG-2 Main Profile@Main Level

© 1998, Philips Research North America, Inc. All Rights Reserved. Philips University

MPEG

MPEG-4 Visual : VOP Encoder



- VOP Encoder is based on each VOP
- Shape Coding, Motion Est.Comp., Content-based Texture Coding Parts

© 1998, Philips Research North America, Inc. All Rights Reserved. Philips University

MPEG

Concept, Development, Mass Production, and Applications of Artificial Retina Chips

Kazuo Kyuma

Mitsubishi Electric Corporation, Japan

Email: kyuma@qua.crl.melco.co.jp

Images of the real world contain a very large amount of informations which present image processing systems cannot analyze in real time with reasonable cost and low power consumption. In order to solve these problems, we have proposed and developed artificial retina chips (AR chip) which combines video camera function (image sensing) and image processing function, in a similar way to the human eyes. The basic principle of our AR chip is based on the novel type of optoelectronic vector/matrix multiplication. The core of the AR chip consists of the two-dimensional array of variable sensitivity photodetection circuit which have also been proposed by us. The matrix corresponding to the input image can be processed in analogue by applying appropriate electric vector to the chip. The possible image processing functions include conventional image sensing, edge extraction, variable resolution, 2D to 1D image compression, pattern matching, noise elimination, and random acces.

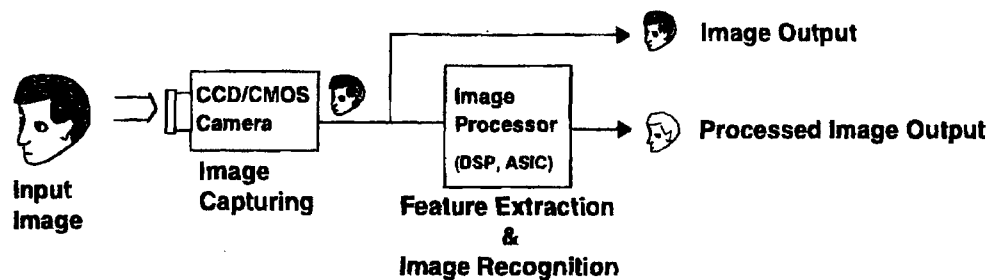
Several types of the AR chips with different resolution have been developed by using the CMOS technology. Among them, we have succeeded in the mass production of 128_5B!_(J128 picels in February 1998. The features of the AR chips over the conventional CCDs are fast image processing (on chip image processing), low power consumption, low cost. With use of these advantages, 500 million chips have been used up to date in several aplication areas which include games, security systems, communication systems, etc. In this talk, the concept, structure, operation principle, some applications, future trends of the AR chips will be introduced.

Concept, Development, Mass Production, and Applications of Artificial Retina chips

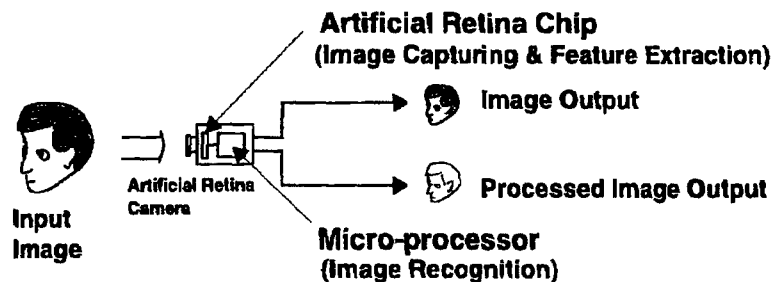
Kazuo Kyuma

**Mitsubishi Electric Corporation
System LSI Division
Advanced Technology R&D Center**

Concept of Artificial Retina Chip

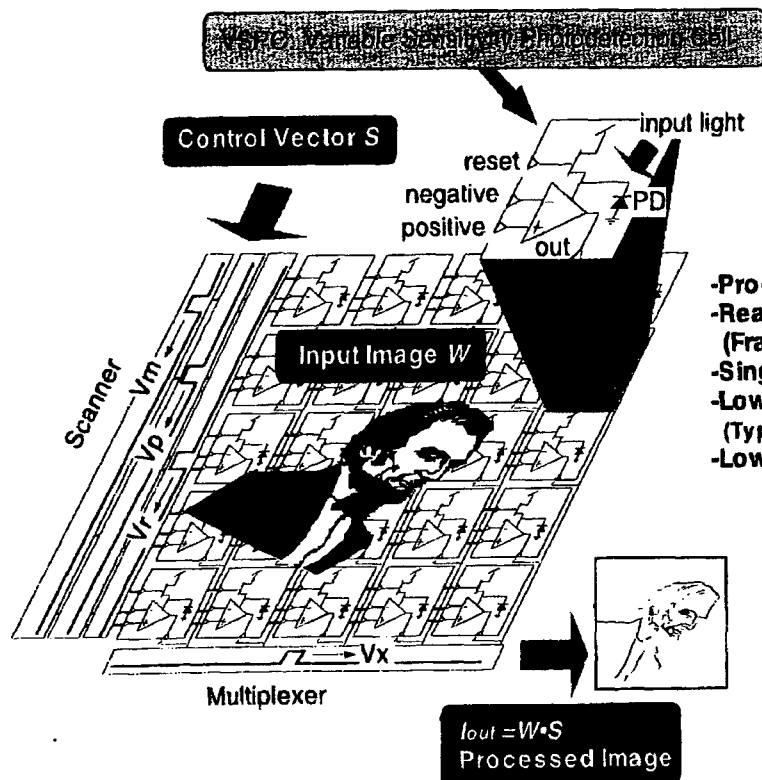


(a) Image Processing System Using CCD or CMOS Sensor



(b) Image Processing System Using Artificial Retina Chip

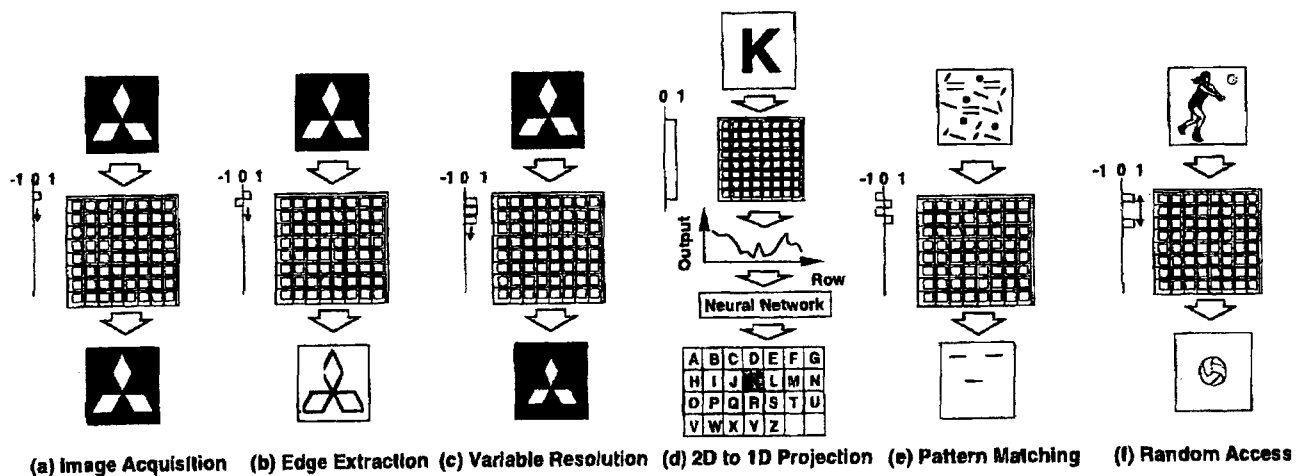
Configuration of Artificial Retina Chip



Features

- Programmable image processor
- Realtime image processing (Frame Rate: 1-500Hz for 32x32 pixel version)
- Single power supply (5V or 3.3V)
- Low power consumption (Typ. 15mW for 32x32 and 128x128 pixel version)
- Low cost

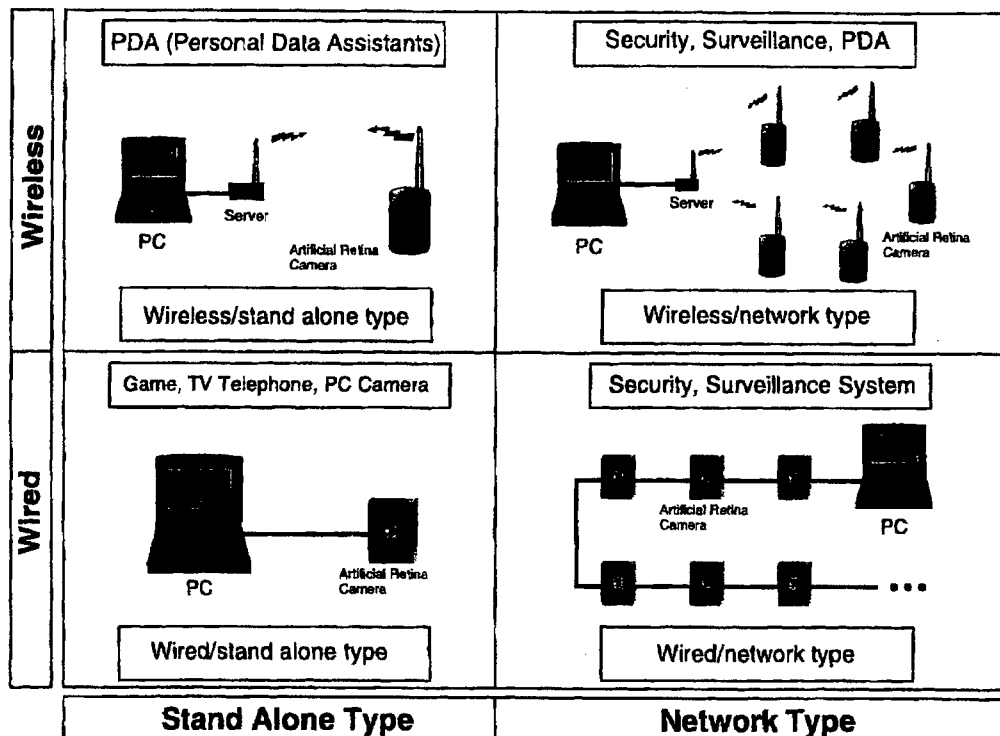
Image Processing Examples



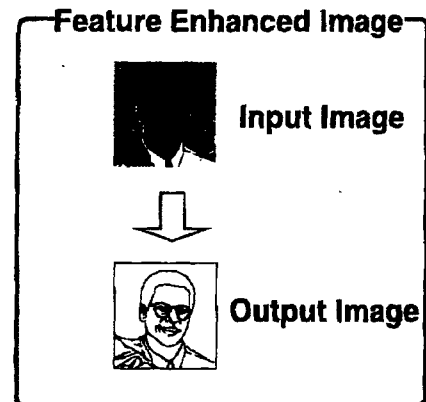
Artificial Retina Chip / Module

Name	Group (Function)	Function outline	PKG	Applications
M64283FP	Artificial Retina Chip	128x128 pixels edge enhancement / detection 2D to 1D projection Random access	16C9-B	Game, PDC PC interface etc.
M64283K	Artificial Retina Chip	128x128 pixels Edge enhancement / detection 2D to 1D projection Random access	20 pins Ceramic SOP	Surveillance camera, security etc.
M64285FP	Artificial Retina Chip	32x32 pixels Edge enhancement / detection 2D to 1D projection Variable data rate	10C2-C	Game, PC interface etc.
M64287U	Artificial Retina Chip	352x288 pixels 2D filtering Built-in AD converter	36 pins Ceramic LCC	Security etc.
M64289U	Artificial Retina Chip	352x288 pixels 2D filtering Built-in AD converter Color type	36 pins Ceramic LCC	Game, PC camera etc.
PCA6050 AG01-01A PCA6100 AG01-01A	Artificial Retina Module	Lens unit, Mitsubishi CMOS Image Sensor, and MCU are installed.	—	Game, PC interface, security etc.
Algorithm	Algorithm for Artificial Retina	Optical flow Pattern / object recognition Numerals / bar code recognition Position recognition	—	Game, PC interface, security etc.

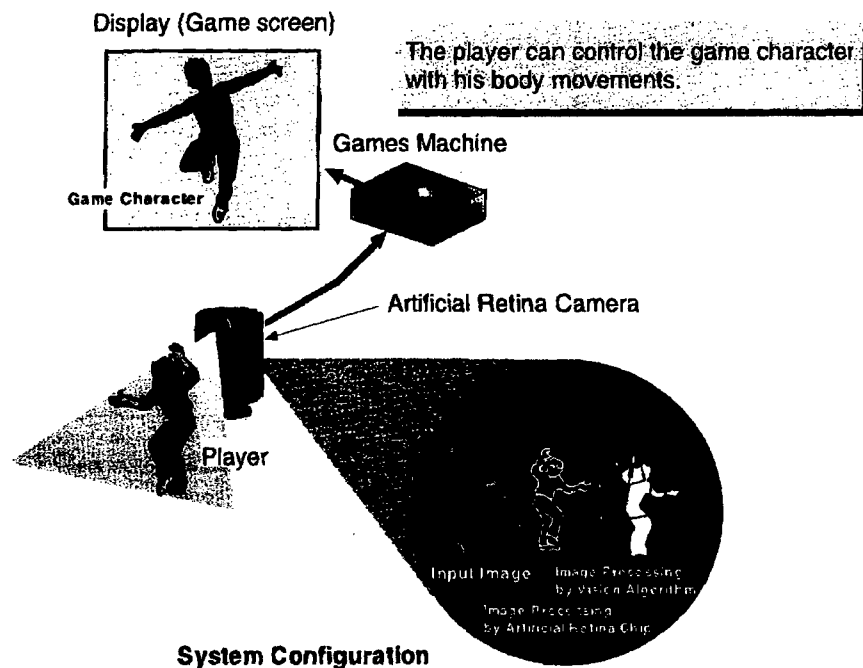
Classification of Artificial Retina Module and Camera



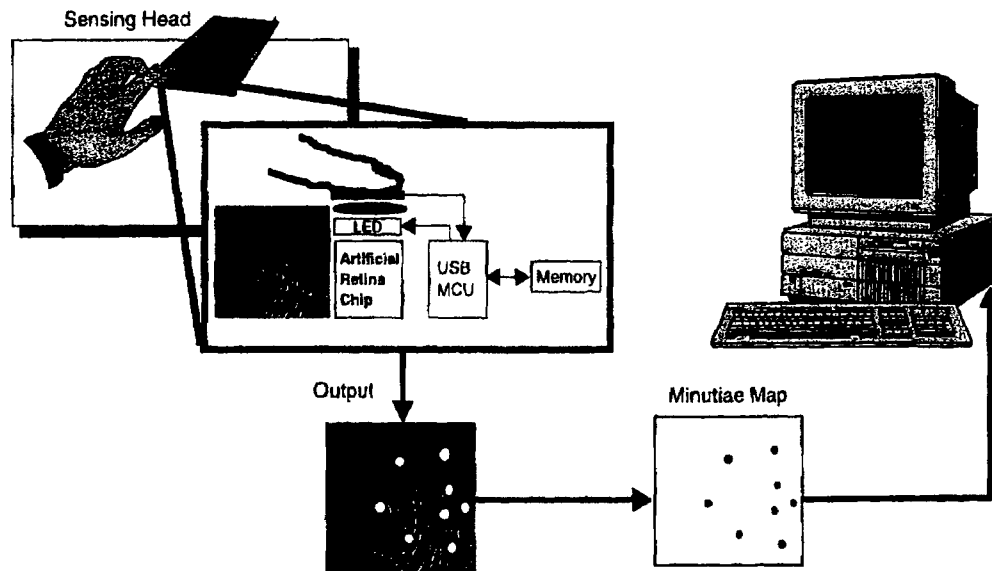
Nintendo Game Boy Pocket Camera



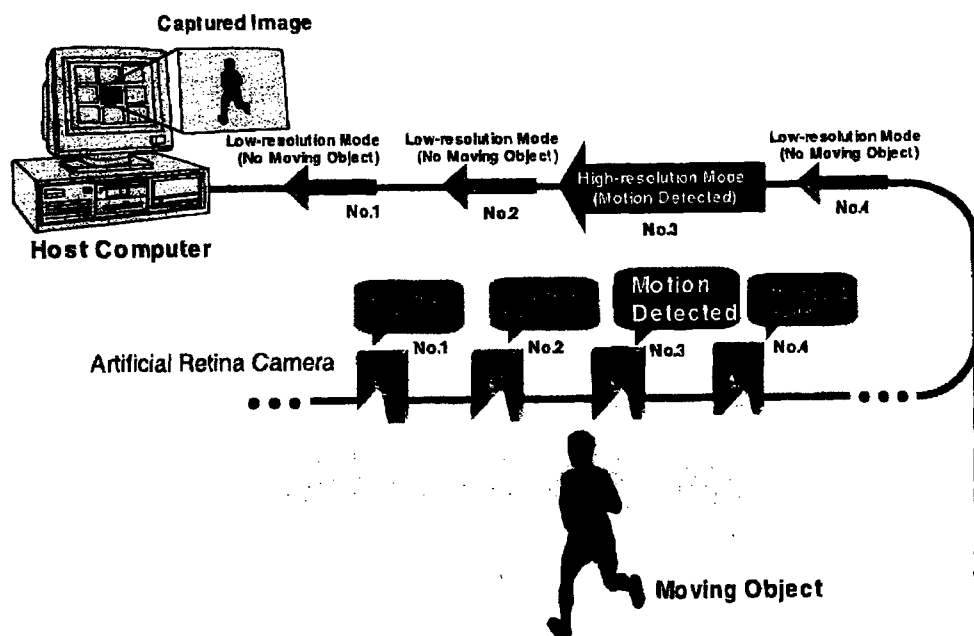
Interactive Game Controlled by Artificial Retina Chip



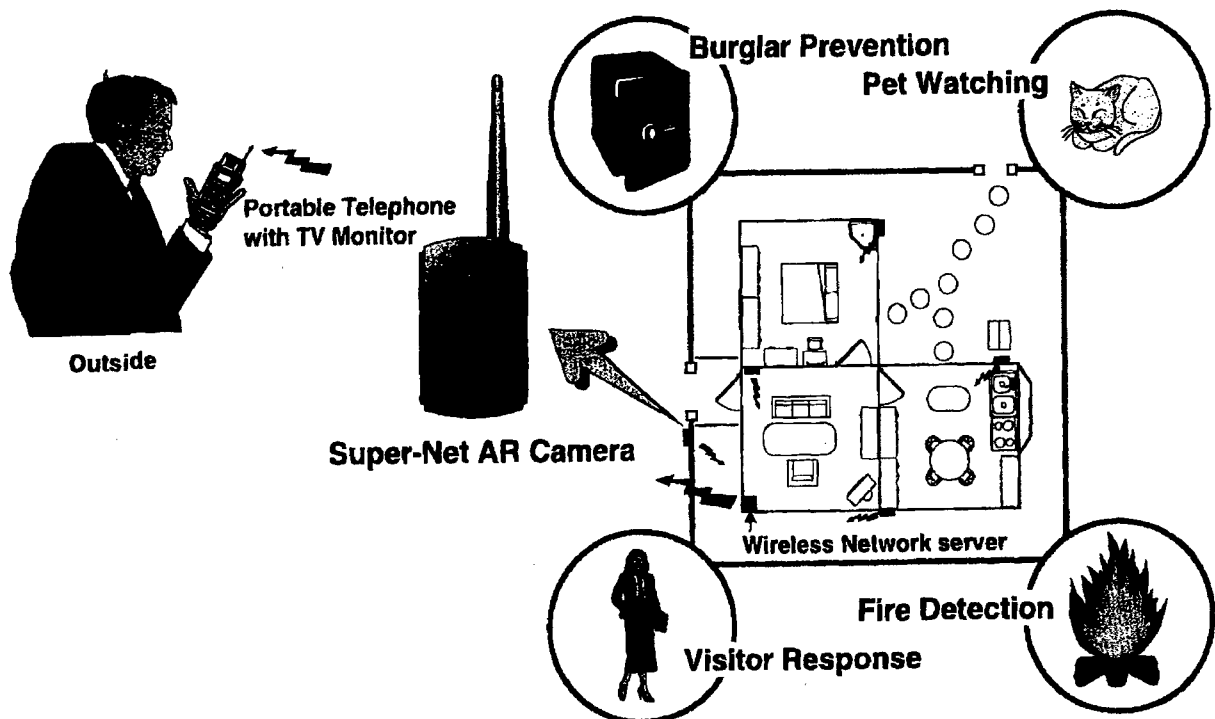
Fingerprint Recognition System



Security Network System



Application of Super-Net AR Camera (Personal Home Security System)



Future Development

1. High resolution (Planning)
VGA:640x480 pixels
2. High speed (Planning)
Frame rate:100Hz
3. Built in USB I/F (Planning)
4. Built in image compression function (Planning)

Data Reduction via Auto-Associative Neural Networks

Claudia V. Kropas Hughes

Air Force Research Laboratory, Materials Directorate,
Wright-Patterson AFB, OH, 45433-7746, USA

Image analysis is a very complex process; many of the relationships are difficult to categorize, much less to program into a computer. The selection of features is the most challenging problem of image analysis, process discovery or sensor fusion. The features must be a data representation that will discriminate the information of interest from the rest of the image. A Neural Network can be a tool for rapid processing of data. Auto-associative Neural Networks (AANNs) are a form of self-organizing maps which can be used to reduce the dimension of the input data in a self-organizing fashion. Dimension reduction is closely related to feature extraction. Features are those datum that efficiently capture the information contained in the entire data set. The data set, has a "superficial" dimensionality of n , and the reduced space of the features that contain all the information about the data has an "intrinsic" dimension of m , where $n > m$. With this property, an AANN can be used to reduce an n -dimensional space to something more intrinsic to the actual data.

Keywords: Image processing, Autoassociative Neural Networks, Feature extraction, Feature selection

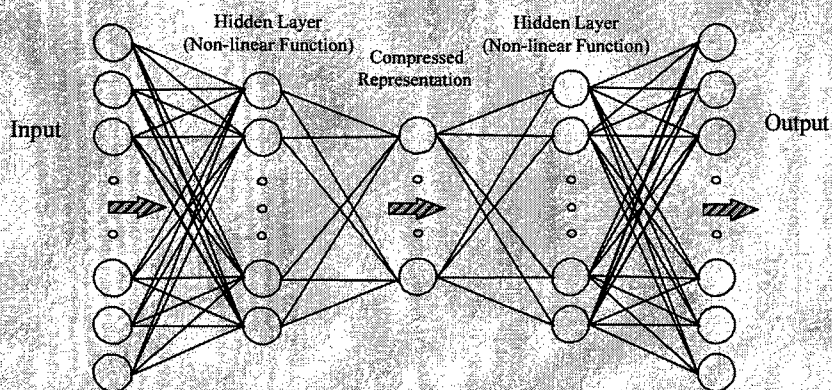
Autoassociative-Heteroassociative Neural Networks For Materials Data Processing

Claudia V. Kropas-Hughes
Air Force Research Laboratory
Materials and Manufacturing Directorate

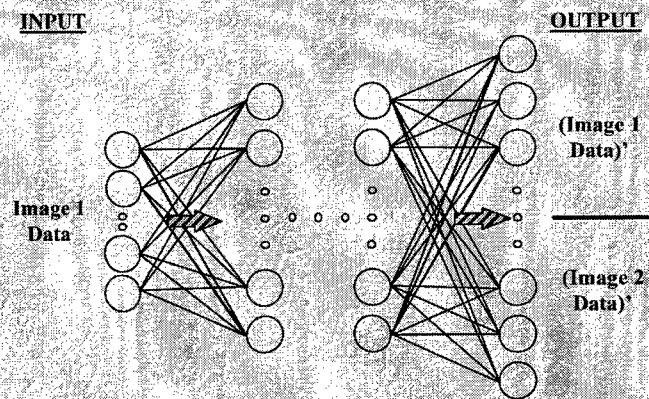
Autoassociative & Heteroassociative Neural Networks (AANN/HANN)

AANN -- $Output = Input$

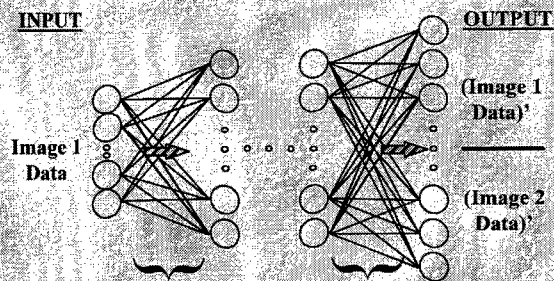
HANN -- $Output \neq Input$



Autoassociative-Heteroassociative Neural Network (A-HNN)



A-HNN



Using 1 Image Data Set for the Encoding Portion of the Network Constrains the Hidden Layer Representations

Constraints from the Encoding Portion of the Network means the Decoding Transformations are Consistent with the Original Data AND the Second Data Set

A-HNN

Two Features of the New Architecture

- Stability of AANN portion used for Network Generalization Robustness check
- Use of input features as target outputs IMPROVES training performance

A-HNN

Stability of AANN

Mathematically, stability of an AANN is determined by the infinity norm metric, that is the distance between the target values and output values achieved by the fully trained AANN. Specifically, for x , the target output value, and z the actual network output, where z is regarded as a weak perturbation of x , the AANN must be constructed such that

$$\text{dist}(x,z) < \epsilon \quad \epsilon > 0$$

For every AANN constructed, in which every point z is locally asymptotically stable, then the AANN is in performance of the training data range, and is interpolating and not extrapolating the information. (Reimann, 1998) This work demonstrates a means to validate the operation of an AANN by evaluating the distance between the target values and network outputs for each feature, for each training sample. An evaluation of the differences, the ϵ values, will determine the stability of the AANN. Any difference values greater than ϵ , will demonstrate that the AANN is operating outside the stability range, e.g. outside the range upon which it was trained by the training data set.

A-HNN: Improved Training

• Improved Training Performance

- Comparisons with conventional Multi-Layer Perceptron (MLP) classification problems
 - Duplication of input features as target outputs provides a network that converges faster
 - fewer training epochs, or
 - fewer flops, or
 - smaller network architecture

A-HNN

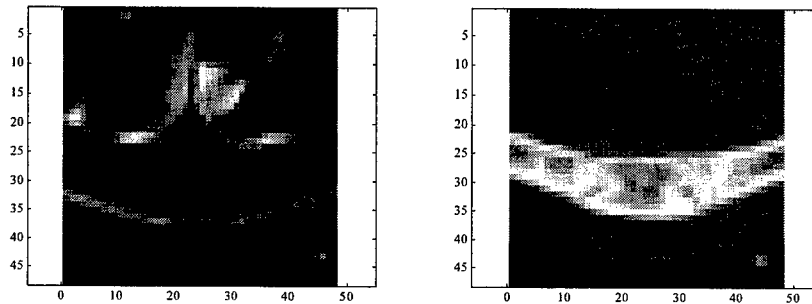
• Three Examples

- Image Processing
- For Classification
 - XOR Data
 - Material Properties Data

A-HNN: Image Processing

1. Feature Extraction: Human Visual System Concepts Applied
 Fourier and Gabor: Limited
 Wavelets: Good compromise
2. Associate the same features by the A-HNN

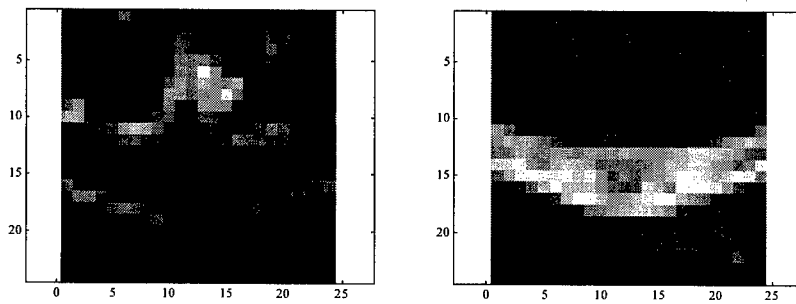
Original images:



Feature Extraction

A-HNN: Image Processing

Low-Low frequency pass wavelet images (1/4 the original size)



Feature Extraction

A-HNN: Image Processing

Line Plots of Col 10 from each of the CT and MR wavelet images:
Note the bone vs soft tissue features are easily separated.

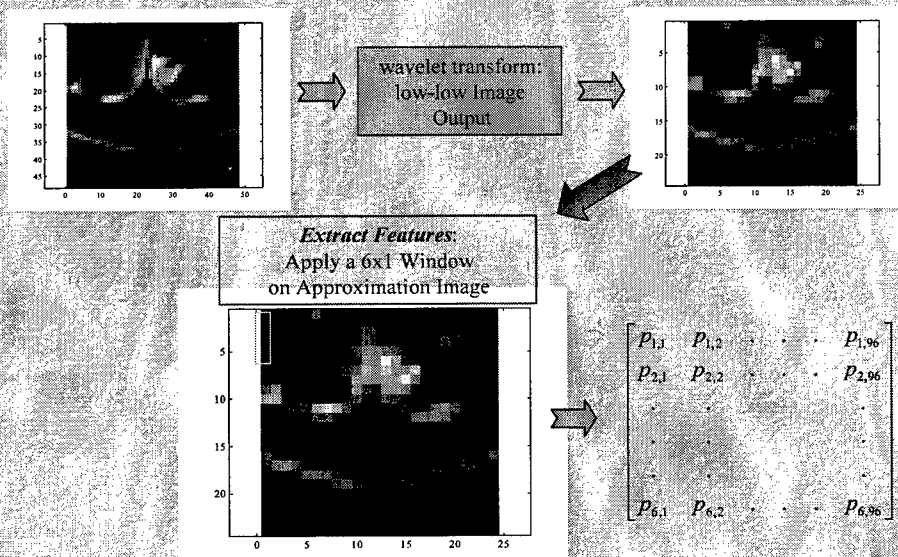
Title:
col10.eps
Creator:
MATLAB, The Mathworks, Inc.
Preview:
This EPS picture was not saved
with a preview included in it.
Comment:
This EPS picture will print to a
PostScript printer, but not to
other types of printers

CT Data

MR Data

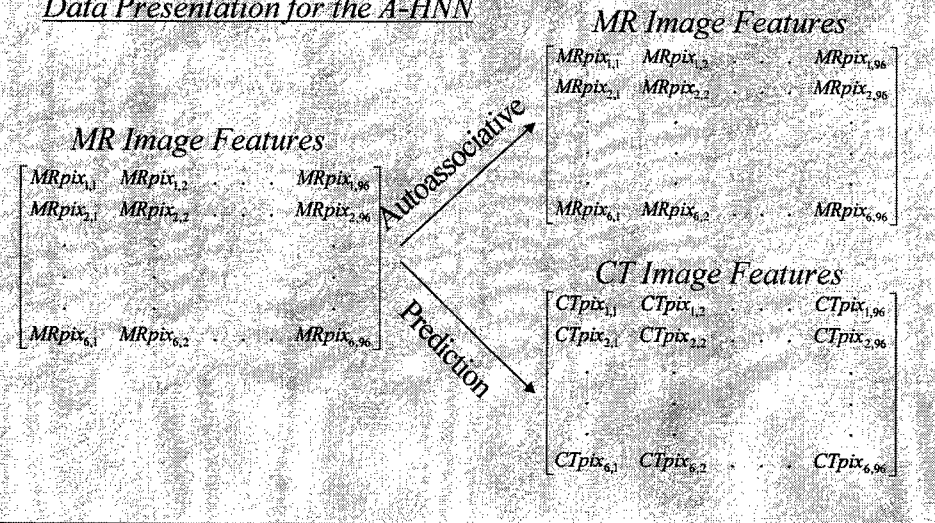
Feature Extraction

A-HNN: Image Processing

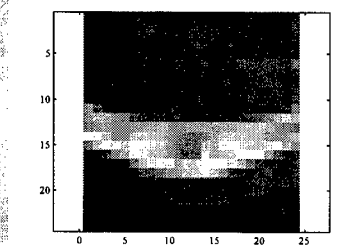
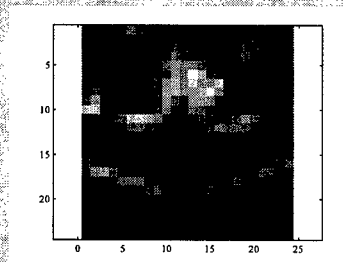


A-HNN: Image Processing

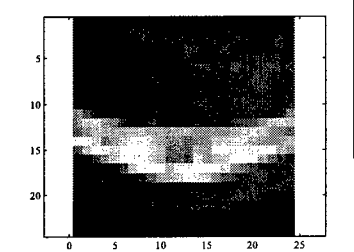
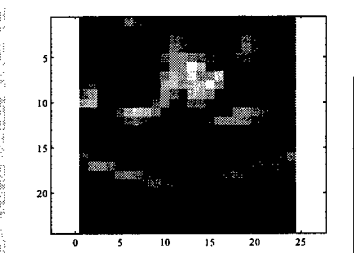
Data Presentation for the A-HNN



A-HNN: Image Processing



A-HNN: Image Processing



Network Output Check:

- *Visual Subjectivity*
- *Distance in Fourier Space*
- *VDP Pixel-by-Pixel Comparison*
- *Network Stability of Autoassociative portion*

A-HNN as a Classifier

• XOR Data Set

• Material Properties Data Set

A-HNN as a Classifier: XOR

Network Output → Test Data Label ↓	Class 0	Class 1	Undecided
Class 0	150	0	0
Class 1	0	150	0

Multi-Layer Perceptron Confusion Matrix for XOR Data using ten nodes on the hidden layer each with sigmoid activation functions.

MLP trained

10 hidden layer nodes
1000 epochs
MSE of 0.320767
78.6 million flops

Network Output → Test Data Label ↓	Class 0	Class 1	Undecided
Class 0	150	0	0
Class 1	0	150	0

A-HNN Confusion Matrix for XOR Data, using the first input value as a target output value and five nodes on the hidden layer each with sigmoid activation functions.

A-HNN trained

5 hidden layer nodes
1000 epochs
MSE of 1.29785
52.1 million flops

A-HNN as a Classifier: Materials Data

- Ternary system (3-element compounds) material property data (Forming and Non-forming Compounds)
- Need to predict the ternary systems that will FORM new compounds
- Current Feature Extraction:
 - Material properties gleaned from journals and books
 - Often Non-reproducible/Non-Verifiable
- Use A-HNN to associate the Material properties to Chemical Elemental Properties (atomic characteristics of each element): Associate the Material Properties with the Physical Nature of the Material
- Translate the experimentally determined features to physical property data that may provide more accuracy in classifying

A-HNN as a Classifier: Materials Data

4104 chemical elemental
properties vectors
15 features in length
3500 used for training

Target values:
15 features in length
-number of electrons in the core, s-, p-,
d-, and f-shells for each element

Input Feature Vector
Chemical Element Properties

a_1	
a_2	
Element(a) a_3	
a_4	
a_5	
b_1	
b_2	Subscript Variable
Element(b) b_3	1 - Number of Valence Electrons
b_4	2 - Energy
b_5	3 - Zunger Ratio
c_1	4 - Melting Point
c_2	5 - Atomic Number
Element(c) c_3	
c_4	
c_5	

Target Output Vector
Atomic Characteristics

a_1	
a_2	
Element(a) a_3	
a_4	
a_5	
b_1	
b_2	Subscript
Element(b) b_3	c - Number of Core Electrons
b_4	s - Number of s-shell Electrons
b_5	p - Number of p-shell Electrons
c_1	d - Number of d-shell Electrons
c_2	f - Number of f-shell Electrons
c_3	
Element(c) c_4	
c_5	
c_6	
c_7	

A-HNN as a Classifier: Materials Data

Input Feature Vector
Chemical Element Properties

a_1	Element(a) a_4
a_2	Element(b) b_3
a_3	b_4
a_4	Element(c) c_1
a_5	c_2
b_1	a_1
b_2	a_2
b_3	a_3
b_4	a_4
b_5	a_5
c_1	b_1
c_2	b_2
c_3	b_3
c_4	b_4
c_5	b_5
c_6	c_1
c_7	c_2
c_8	c_3
c_9	c_4
c_{10}	c_5

Inputs: All 15 features
Targets: Reduced subset
of inputs

Input Feature Vector
Chemical Element Properties

a_1	Element(a) a_4
a_2	Element(b) b_3
a_3	b_4
a_4	Element(c) c_1
a_5	c_2
b_1	a_1
b_2	a_2
b_3	a_3
b_4	a_4
b_5	a_5
c_1	b_1
c_2	b_2
c_3	b_3
c_4	b_4
c_5	b_5
c_6	c_1
c_7	c_2
c_8	c_3
c_9	c_4
c_{10}	c_5

Subscripts
1 - Number of Valence Electrons
2 - Energy
3 - Zunger Ratio
4 - Melting Point
5 - Atomic Number
c - Number of Core Electrons
s - Number of s-shell Electrons
p - Number of p-shell Electrons
d - Number of d-shell Electrons
f - Number of f-shell Electrons

Inputs: Reduced
subset of the 15
features
Targets: "Rejected"
input features.

Single-hidden-layer networks with 10 nodes

A-HNN as a Classifier: Materials Data

Network Configuration Comparisons

Network Configuration	Epochs	Hidden Layer Nodes	Network Accuracy	Least Feature Accuracy	Number of Flops
A-HNN	15,000	25	93.6%	83.4%	4.7955 E10
Multi-Layer Perceptron	35,000	15	92.2%	68%	4.3985 E10
A-HNN: Subset of Inputs to Different Subset of Inputs as Targets	30,000	20	88.8%	73.8%	5.3568 E10
A-HNN: Full Set of Inputs to Subset of Inputs as Targets	30,000	10	84.4%	67.4%	3.38619 E10

***A-HNN Best Accuracy
Smaller Network
Fewer Training Epochs***

Publications

• Patent Application

- Autoassociative-Heteroassociative Neural Network

• Ph.D. Dissertation, Air Force Institute of Technology, May 1999

Image Processing Plume Fluence for Superconducting Thin-Film Depositions

J.G. Jones*, R.R. Biggers*, J.D. Busbee*¹, D.V. Dempsey*²,
G. Kozlowski**

* Air Force Research Laboratory, Materials Directorate
Wright-Patterson AFB, OH, 45433-7746

** Air Force Research Laboratory, PRP, Wright-Patterson AFB

¹Technical Management Concepts, Inc. Beaver Creek, OH 45434

²University of Dayton Research Institute Dayton, OH 45409

Process control is a crucial element in all deposition techniques. It is especially elusive in the versatile and efficient deposition technique known as pulsed-laser-deposition (PLD). Image processed emissions from the plume of a $\text{YBa}_2\text{Cu}_3\text{O}_{7-x}$ (YBCO) target are monitored *in situ* to determine two dimensional spatial information about the plume. Manual and fuzzy-logic based regulation of laser energy based on this plume emission feedback resulted in improved film quality and repeatability of the PLD thin-film depositions. Imaging of the plume under various deposition conditions, both with and without process control, will help to improve understanding of the effect of changing environmental conditions on the plume characteristics.

Keywords: Process control, Fuzzy logic, Image Processing, Control non-linearities, Process identification

Image Processing Plume Fluence For Superconducting Thin-Film Depositions

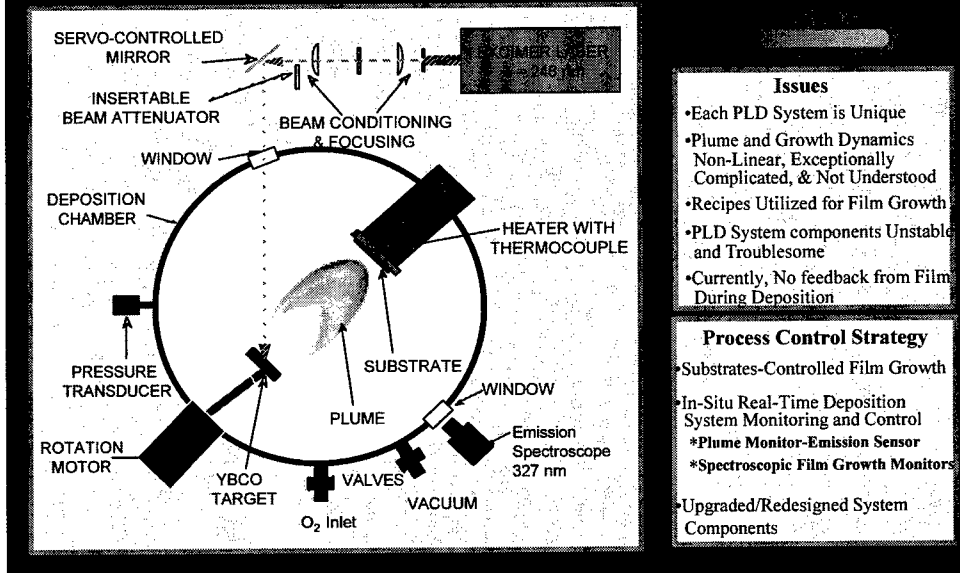
J. Jones, R. Biggers, J. Busbee,
D. Dempsey, G. Kozlowski

Air Force Research Laboratory,
Materials & Manufacturing Directorate

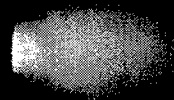
Topics Covered

- Introduction
- Pulsed Laser Deposition (PLD) Process
- *In Situ* Imaging
- Fuzzy Logic *In Situ* Control

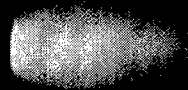
PLD Apparatus



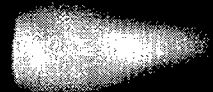
Pressure Effect on YBCO Plume



150 mTorr, 760 C

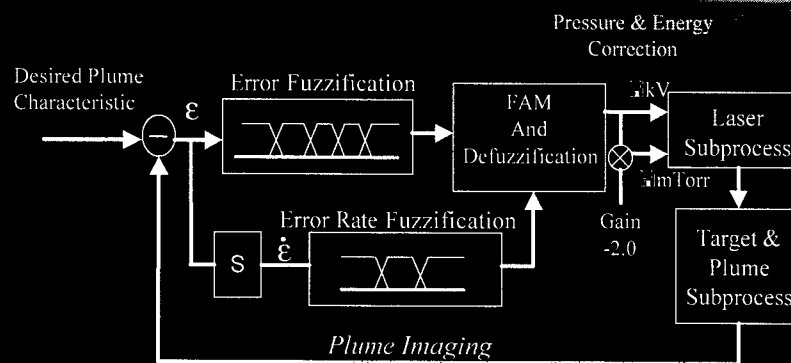


400 mTorr, 760 C



800 mTorr, 760 C

Process Control Based on Imaging



Fuzzy Controller Block Diagram For Spatial Plume Regulation

Method of Computing Actuation Change in PLD Process

Fuzzy Associative Memory Bank For Spatial Regulation

Error Rate	P	PF	PH	NH	NH	NF
	Z	PF	PH	ZO	NH	NF
	N	PF	PH	PH	NH	NF
		NL	NS	ZE	PS	PL
		Error				

$$O = \frac{\sum_{i=1}^{NRules} Z_i \min(X_i, Y_i)}{\sum_{i=1}^{NRules} \min(X_i, Y_i)}$$

Innovations in Materials Design

(abstracts and viewgraphs)

Toward the Future: Innovations in Materials Design

Shuichi Iwata

RACE, Faculty of Engineering, The University of Tokyo, 7-3-1 Hongo, Tokyo 113, Japan

With the advent of personal computing and the internet, the opportunities for realizing new ideas and innovations are virtually limitless. In the context of materials development, an area which is pervasive and impacts nearly all aspects of our lives, these opportunities are potentially very profound. Materials design, as a profession and a field of research, has had an enduring legacy of contributions, involving methods ranging from first principle calculations to mesoscopic and macroscopic techniques, and more recently empirical or data-driven approaches.

Innovation is often associated with a 'leap' in performance or breakthrough from one equilibrium point to another new point with new values. Innovations require 'something' new as well as incremental efforts based upon a feasible and rational plan. It is a dynamic process, driven by a 'sixth' sense and a human desire to discover. This human aspiration to innovate will be driven by and will leverage an 'exploding' information technology age which will foster evermore creative environments for design.

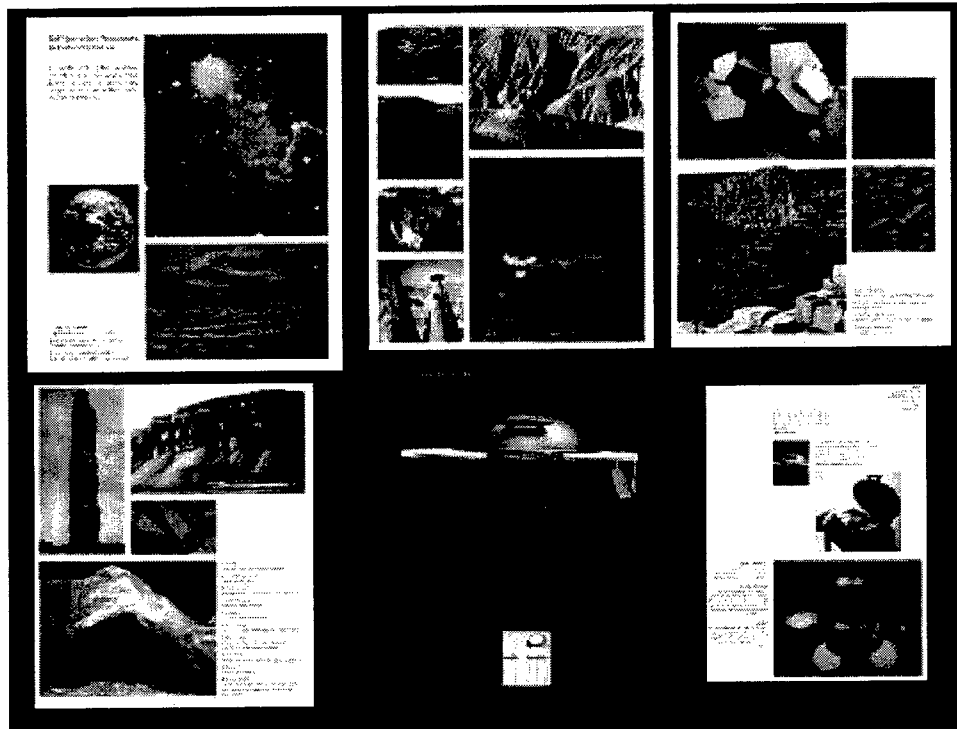
Users of these environments will use tools which evolve and adapt to the user community. The speed and ease of use of such environments will accelerate the fusion and synergy of varying opinions and ideas, while minimizing expensive trial and error. New insights, however unorthodox, may be explored more rapidly, and at minimal expense, via virtual design environments which enable the simulation and modeling of artifacts regardless of their complexity. In this context, we have organized this 'Innovations in Materials Design' session to stimulate your awareness of future directions for materials design, and more generally, we want you to reconsider the 'richness' of the computing resources, data, knowledge, methods, and communication technology available. .

Innovations in Materials Design

RACE, The University of Tokyo
Shuichi IWATA

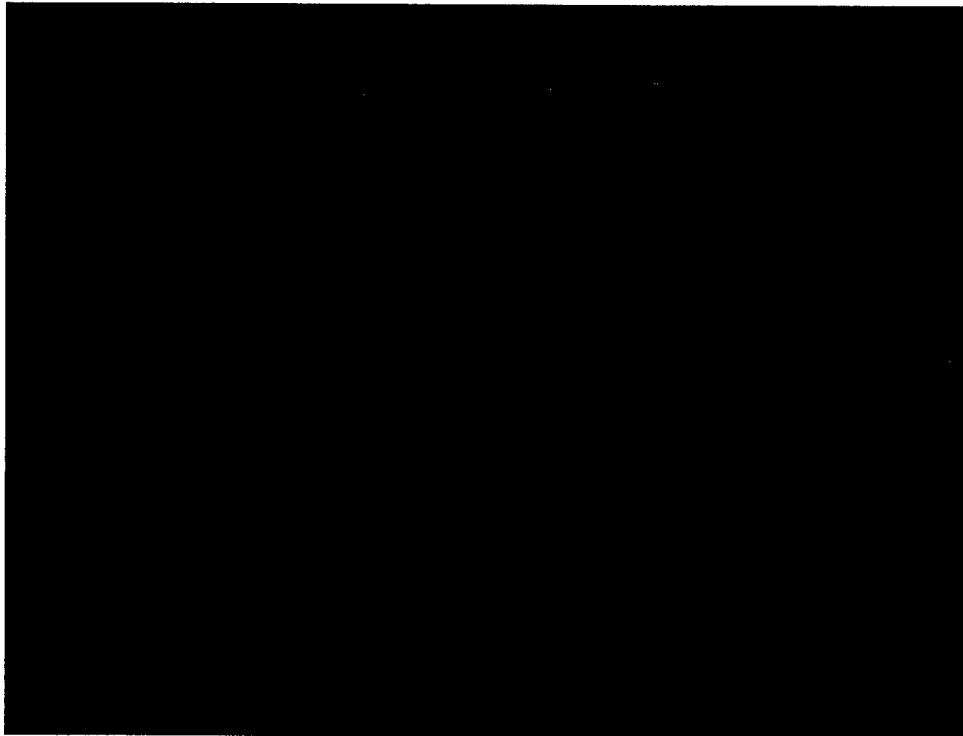
Facts, facts, facts,

Nothing to say... ..

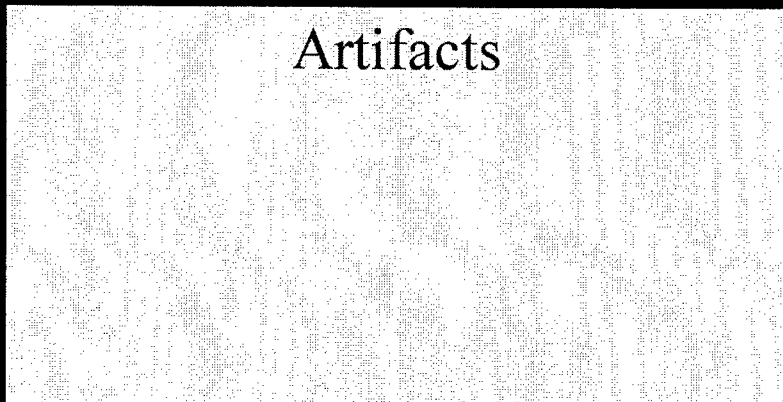


Explanations based on Some Sciences

- Knowledge Primitives
- Constraints
 - Axioms, Theories, Rules,
 - Facts
- Self-organization Mechanism
 - Emotions, Evolutions, Emergence,
- Tentative Explanations to form New Knowledge Primitives



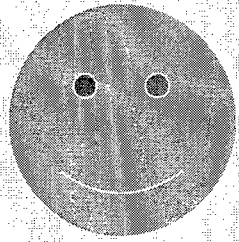
Great Possibilities in Virtual Space



Facts already realized

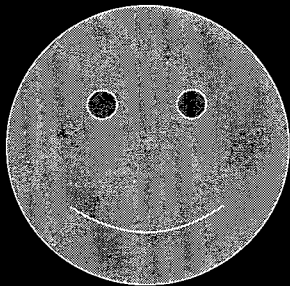
Great Possibilities in Virtual Space

Artifacts

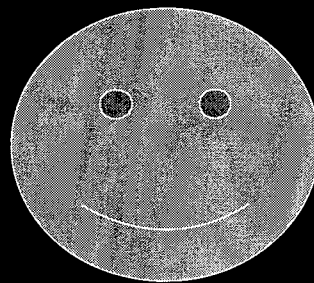


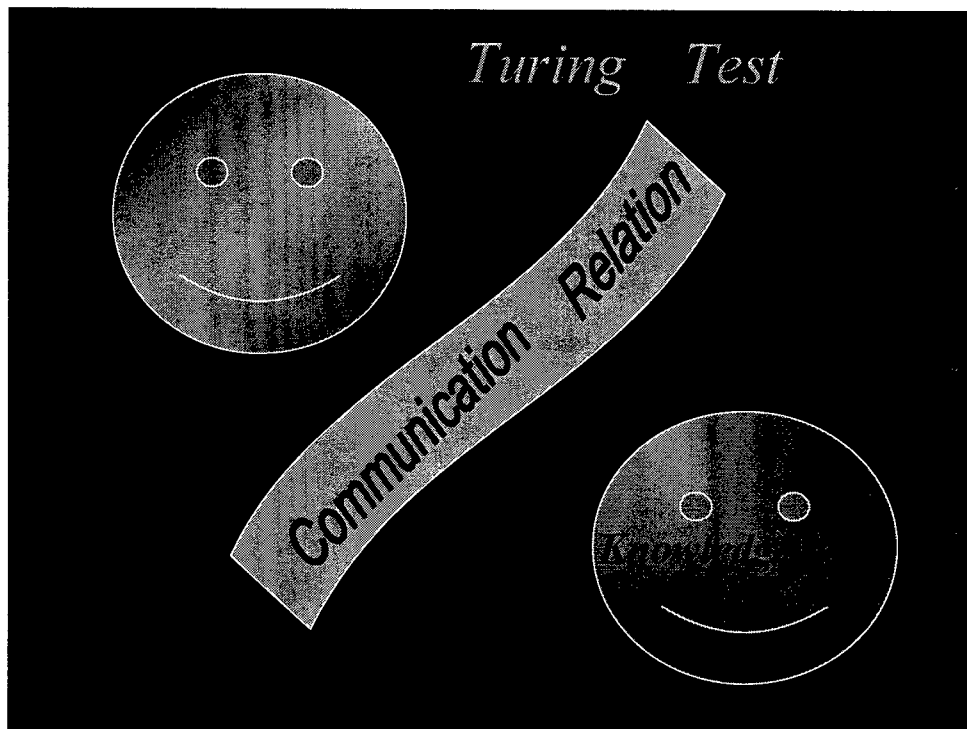
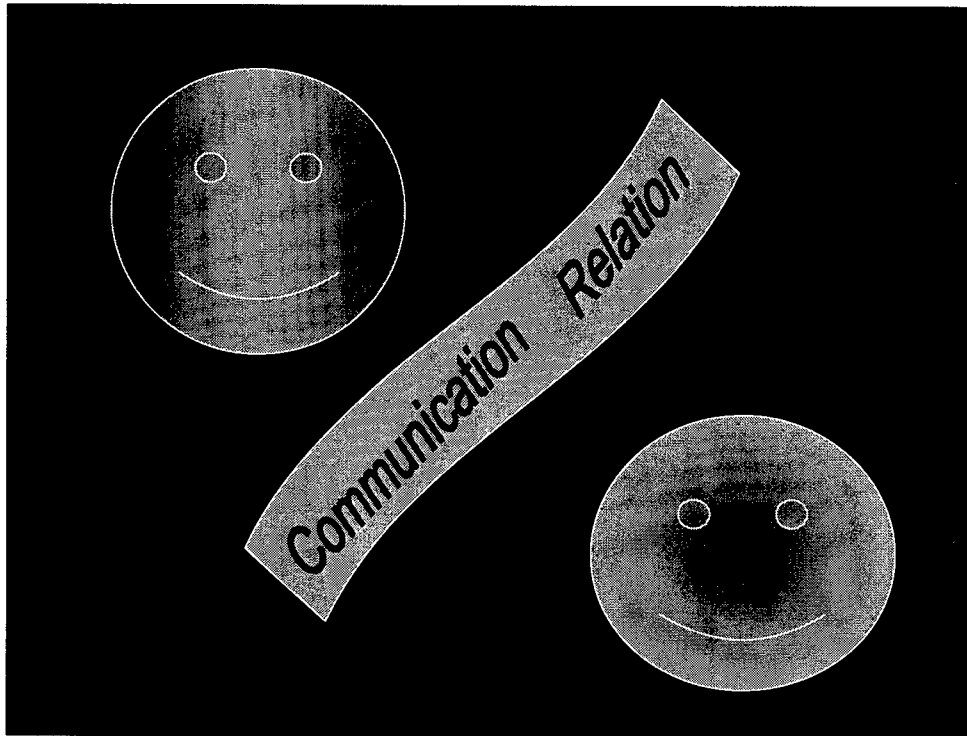
Role of Intelligence

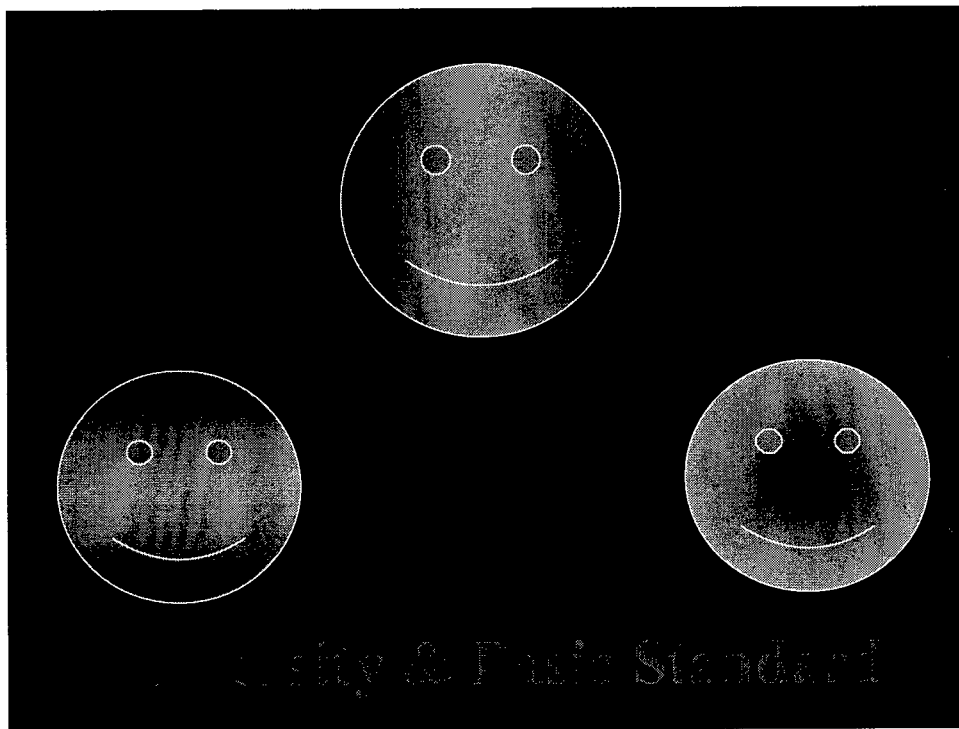
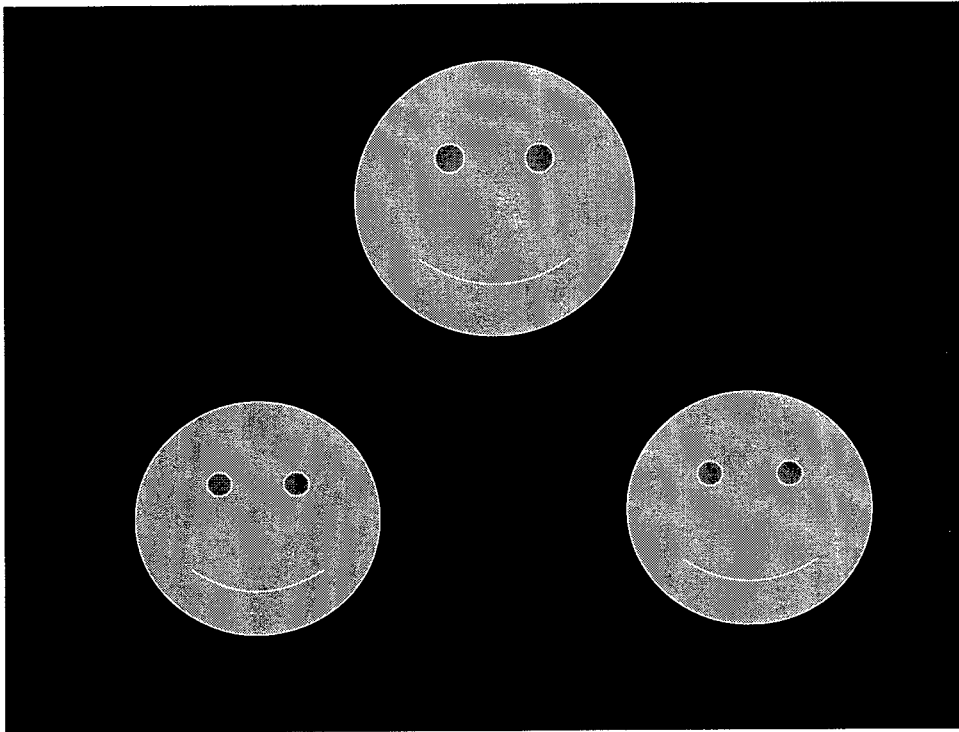
Facts already realized

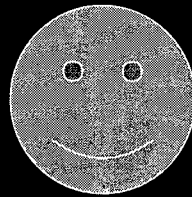
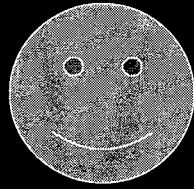


?

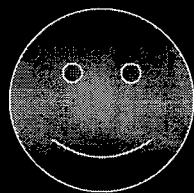
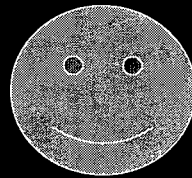
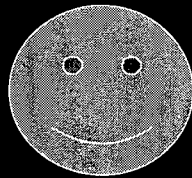








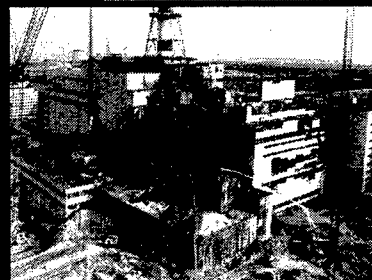
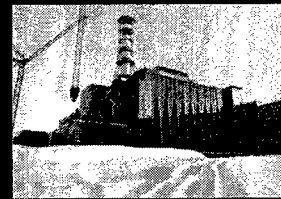
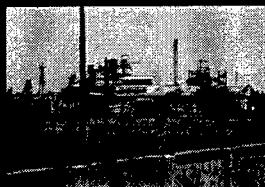
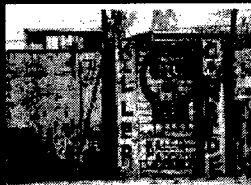
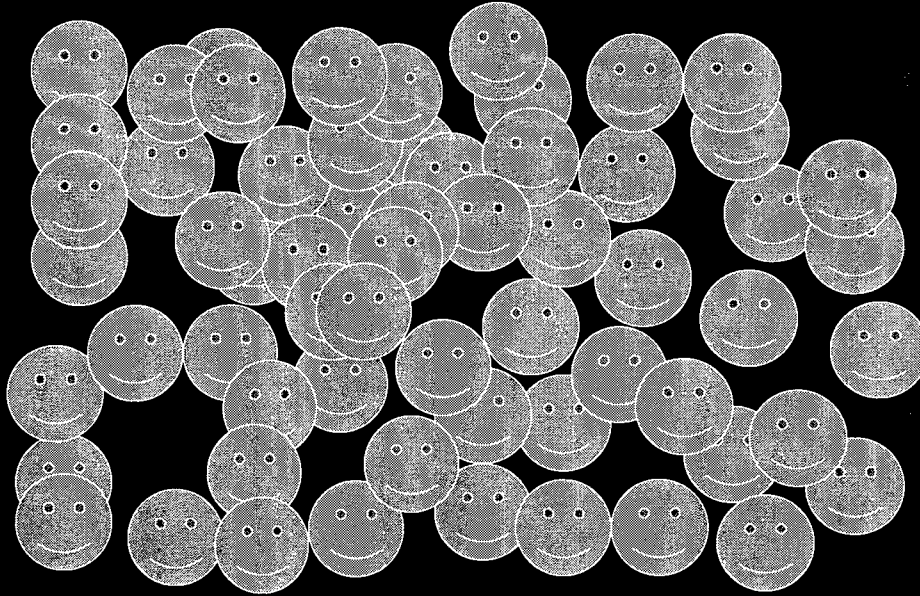
Four Girls in a Dormitory

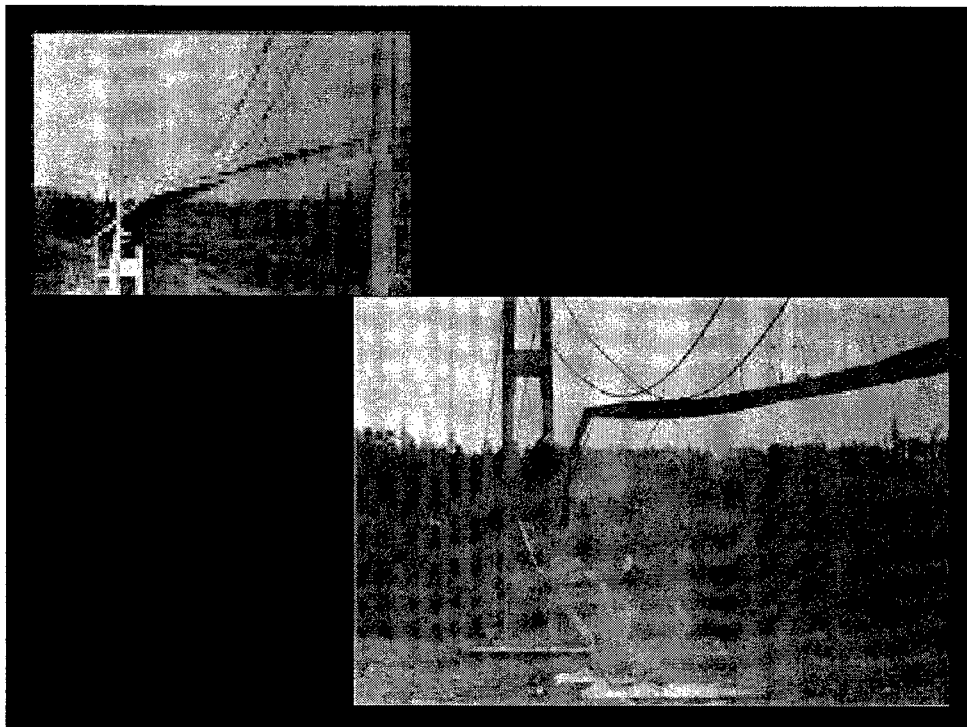
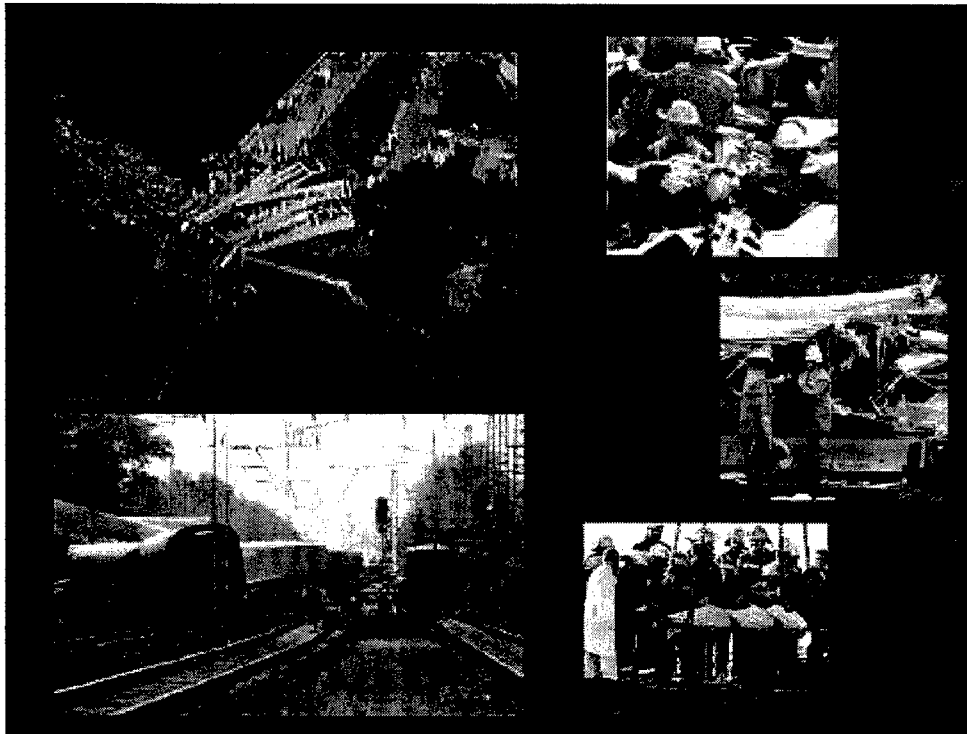


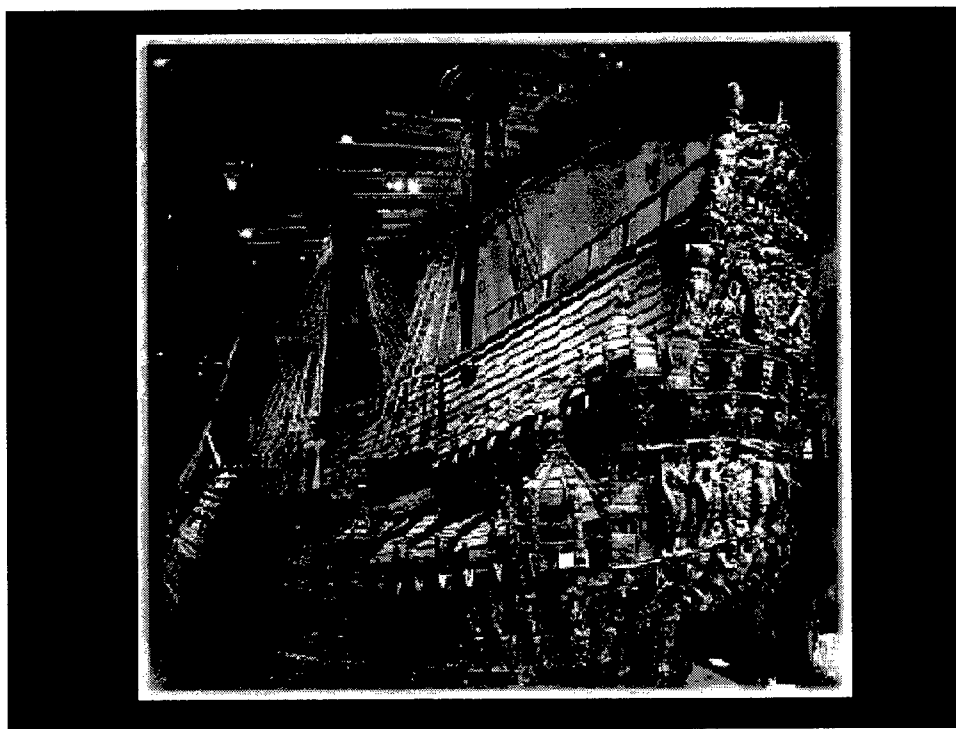
Four Girls in a Dormitory



Issues on Homogenization

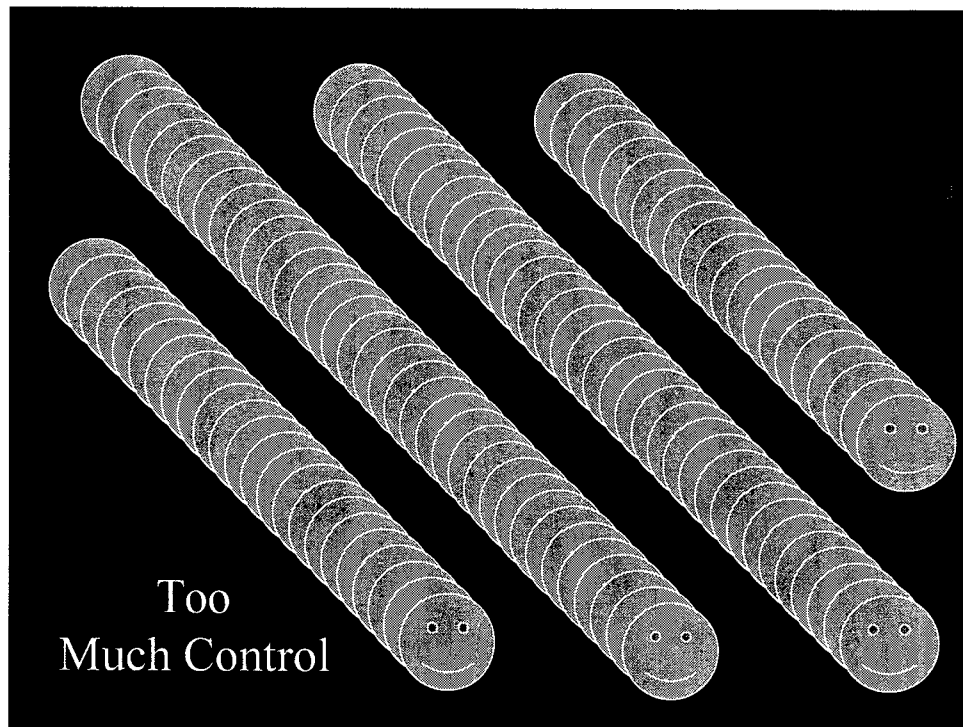




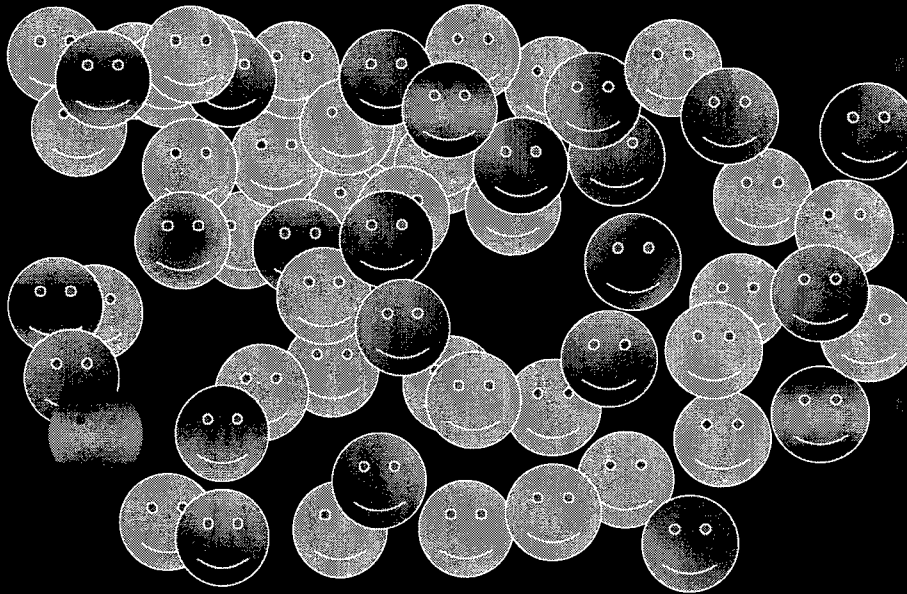


Dangers in WTA Market (winner-take-all market)

- The Productive Edge?
- Total Creativity in Business & Industry?



Dynamism of Group and Creativity



Denkenexperiment

- Di-, Tri-, Quadri-,... ... Lemma
- Paradox
- Over-simplification e.g., Toxicity, Witch Hunting
- Ambiguity Issues
 - Nothing is better than my wife.
 - A penny is better than nothing.
 - Hence a penny is better than my wife.
- Closed Perfect / Open Fluctuation
 - convergence, divergence,

SUN TZU ON THE ART OF WAR
THE OLDEST MILITARY TREATISE IN THE WORLD

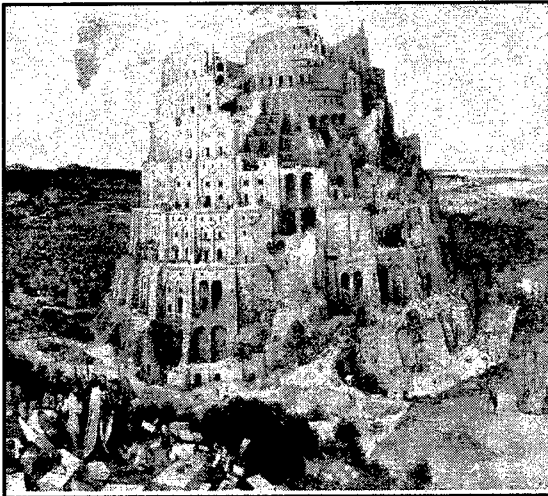
Translated from the Chinese, By LIONEL GILES, M.A. (1910)

I. LAYING PLANS

1... ..,25... ..

26. Now the general who wins a battle makes many calculations in his temple ere the battle is fought. The general who loses a battle makes but few calculations beforehand. Thus do many calculations lead to victory, and few calculations to defeat: how much more no calculation at all! It is by attention to this point that I can foresee who is likely to win or lose.





Atomic Environments in Relation to Compound Prediction

Jo Daams * and Pierre Villars **

* Philips Research, The Netherlands

** Material Phases Data System (MPDS), CH-6354 Vitznau, Switzerland

Predicting new materials and their respective physical properties is the most challenging objective for every scientist working in materials science. Experience, intuition and cooperation with other scientist have been the usual tools of the experimentalist looking for new specific materials. To date, our search for new rules, generally learned from experience, are mostly intuitive and typically undocumented. These unreported empirical approaches have prevented newcomers in material science from making use of these rules and/or regularities.

Our main objective in this study was to collect together as many as possible of these rules or regularities by analyzing all the available published structural data. Combining these rules and regularities with theoretical, practical or empirical models could lead to the desired objective or should at least limit the "scientific area" to where a solution can be found. The first analysis method applied to the published structural data is a purely geometrical analysis. In this study, we analyzed the crystal structure of all intermetallic structure prototypes for their geometrical correctness and, by doing so, we also determined for each atom in the asymmetrical unit - the atomic environment or coordination polyhedron. After completion of this analysis, we were able to define a limited number of "most-frequently-occurring" Atomic Environment Types (AET's).

It was then relatively easy to combine crystal structure data and AETs into coordination prototypes which resulted in a large decrease in the number of structure prototypes. The resulting structure prototypes combine varying crystal structure types into the same AETs and therefore into the same coordination prototype. We also separate out the odd partly-incorrect structure types from the correct ones. We use atomic properties, such as the electron negativity, the number of valence atoms and Zunger's pseudo potential radius for the construction of so-called Structure Stability Diagrams. In these SSD's it is possible to define very sharp three-dimensional volumes wherein compounds have the same crystal structure and, to a certain extend, the same physical properties.

In Quantum Structure Diagrams (QSD's) the observed AETs and the calculated SSD volumes are combined into a 3-D space model which may result in a "semi empirical" prediction method. Although our method is by its nature limited to already observed materials and their physical properties because it is based on published data, studying it thoroughly will most certainly give us more insight into known materials and lead to new, so far unknown, roads to success.

In our presentation we will introduce how the AET analysis is done and present some results for the intermetallic compounds. We will show how SSD's are constructed, using relatively simple mathematical expressions, from the above-named atomic properties and present some results. We will show how these methods are combined into QSD's and present some QSD's which were used to predict new materials.

General Reference

Chapters 11&15 of "Intermetallic Compounds, Principles and Practice", edited by J.H. Westbrook and R.L. Fleischer. Published by John Wiley & Sons (1995) . ISBN 0-471-94219 7.

Atomic environments in relation to compound prediction

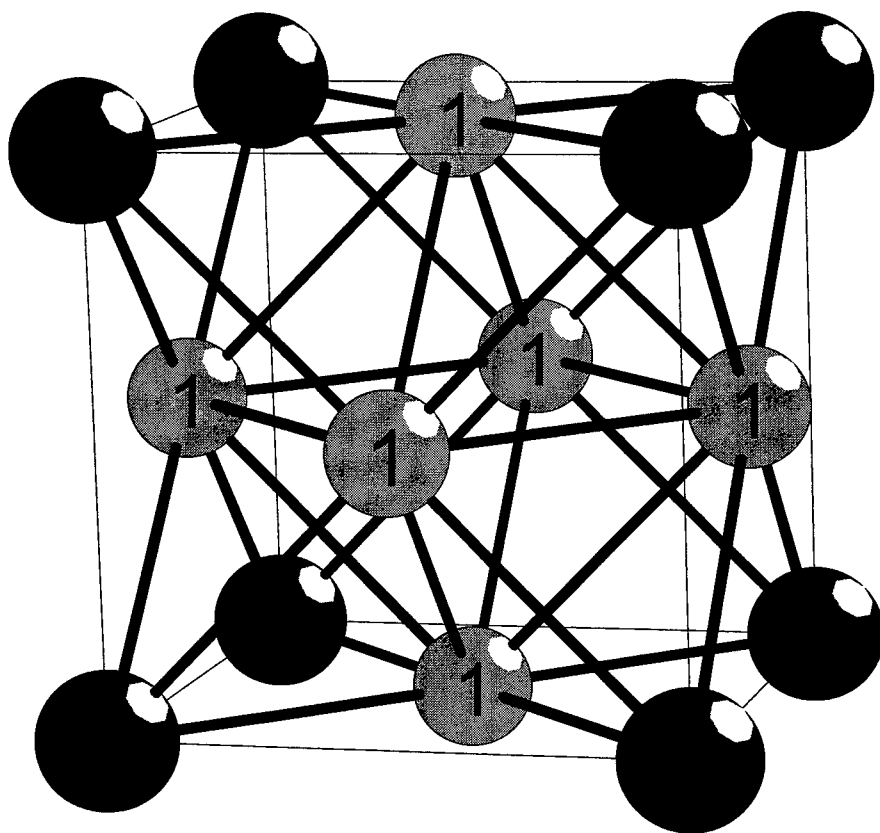
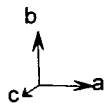
Jo Daams¹ & Pierre Villars²

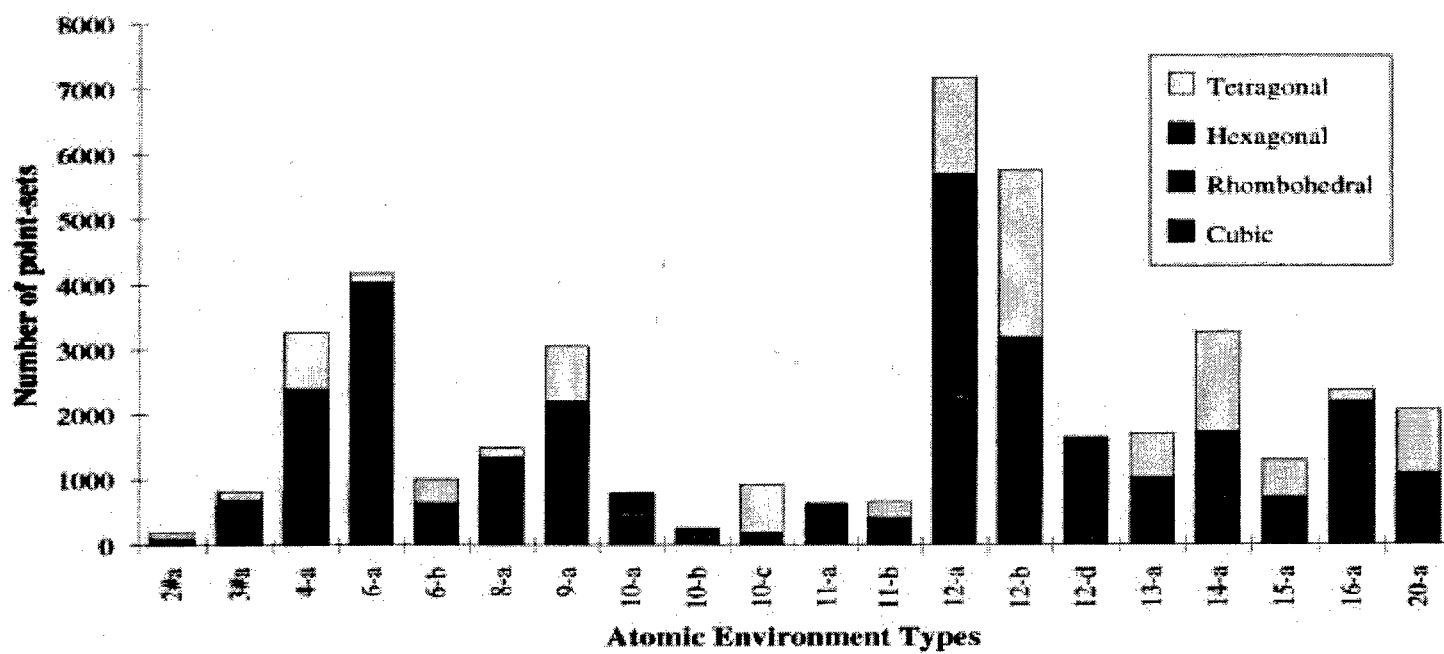
**1. Philips CFT, Prof. Holstlaan 4, 5656 AA
Eindhoven, The Netherlands.**

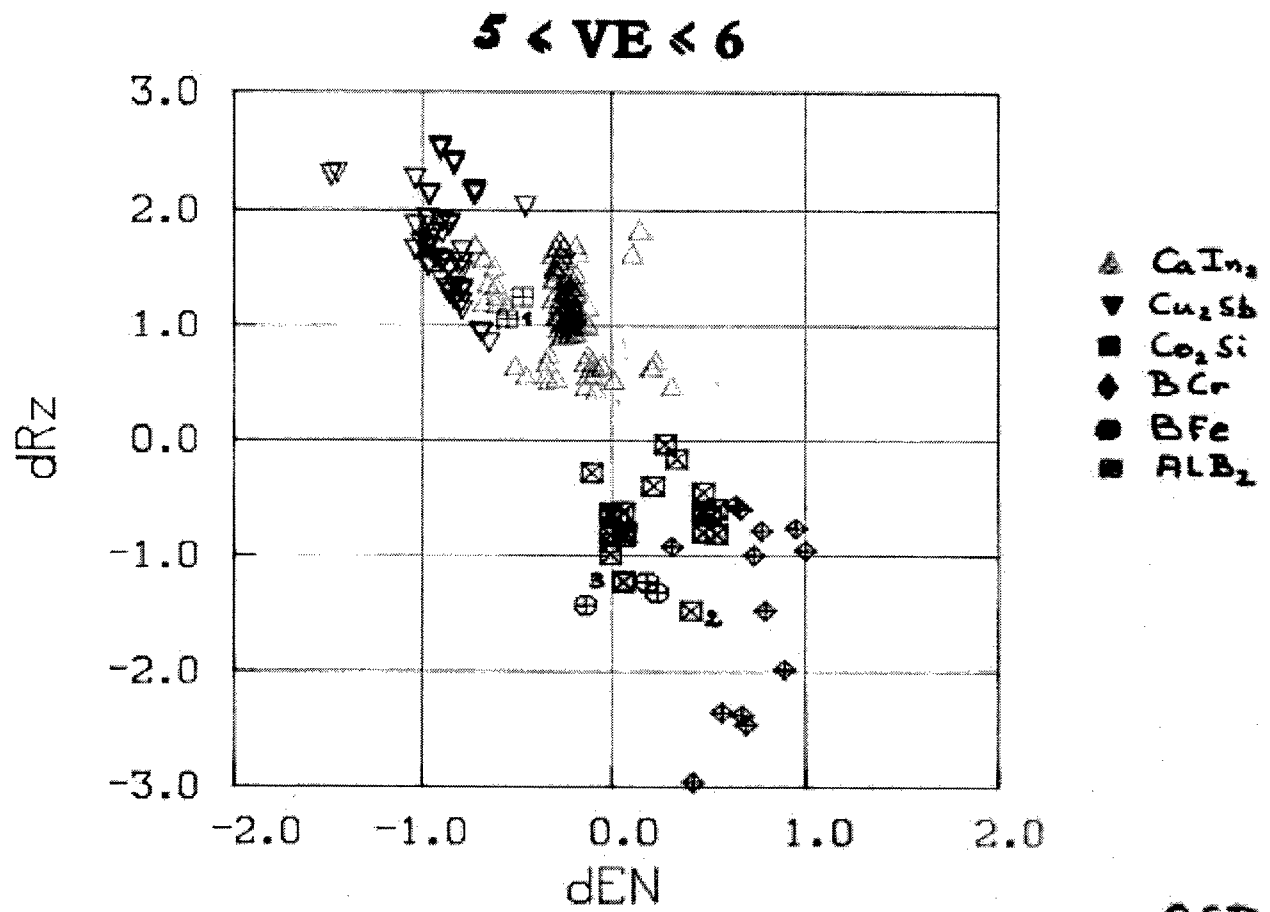
**2. MPDS, Schwanden 400, CH-6354
Vitznau, Switzerland.**

Outline

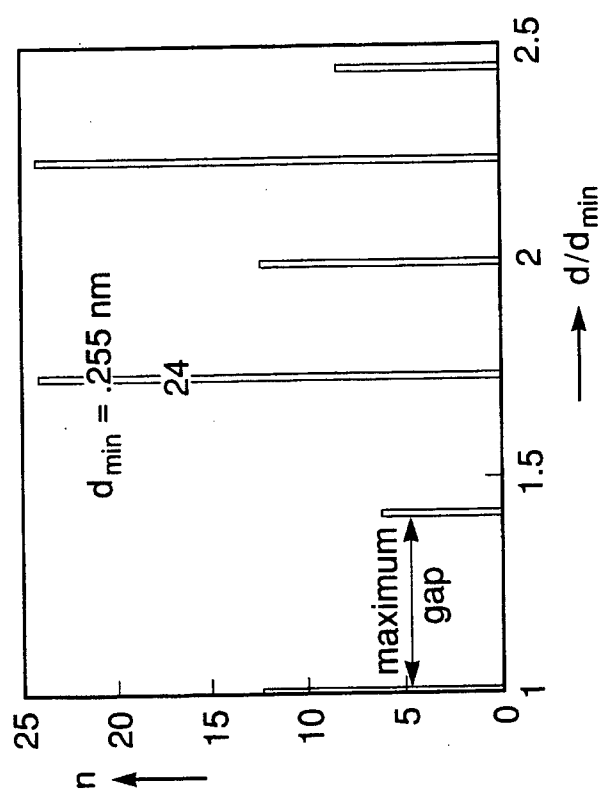
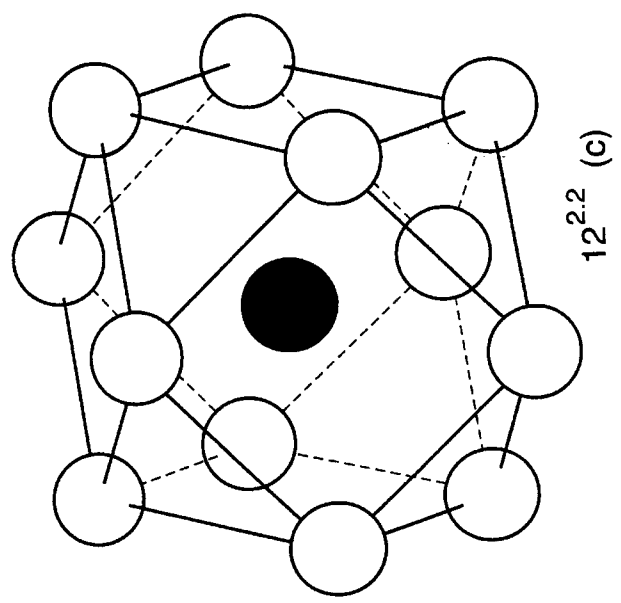
- 1) Atomic environment and co-ordination types**
- 2) Structure stability diagrams**
- 3) Quantum structure diagrams**
- 4) Conclusions**

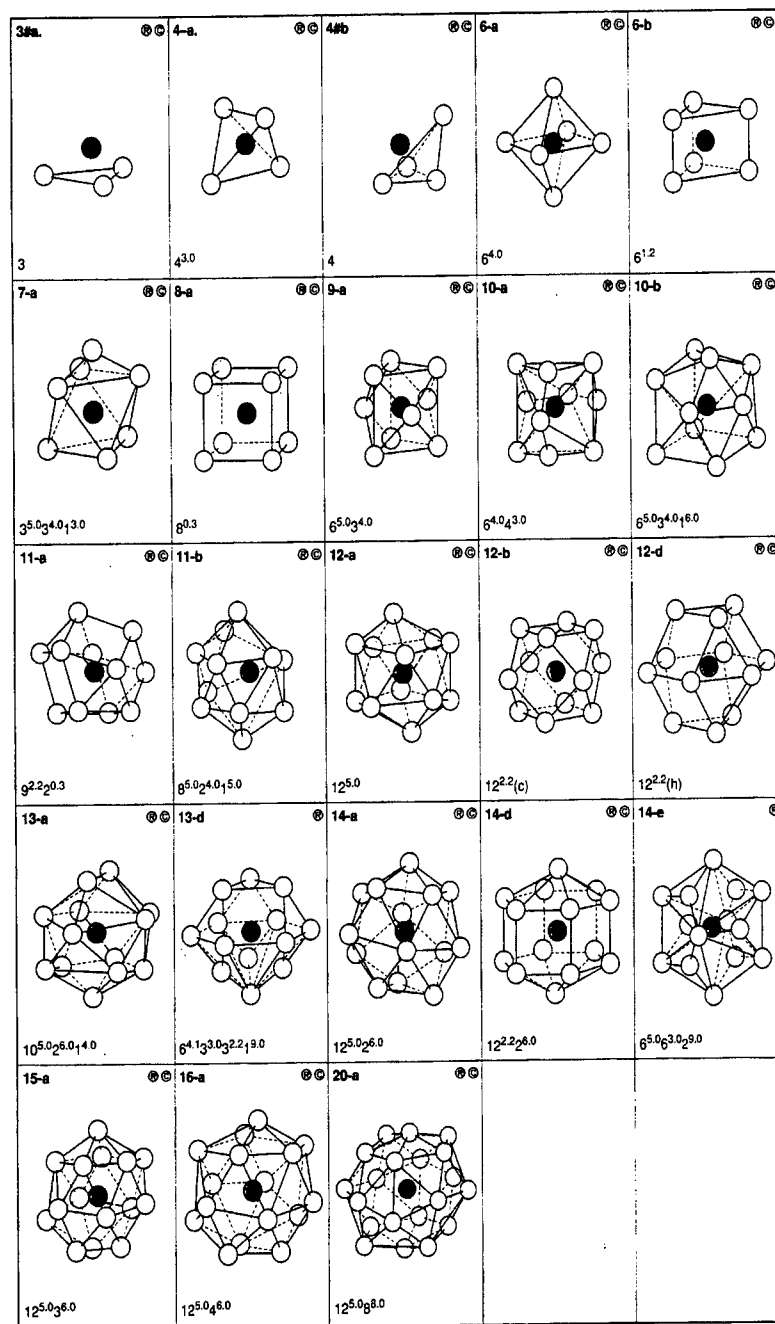


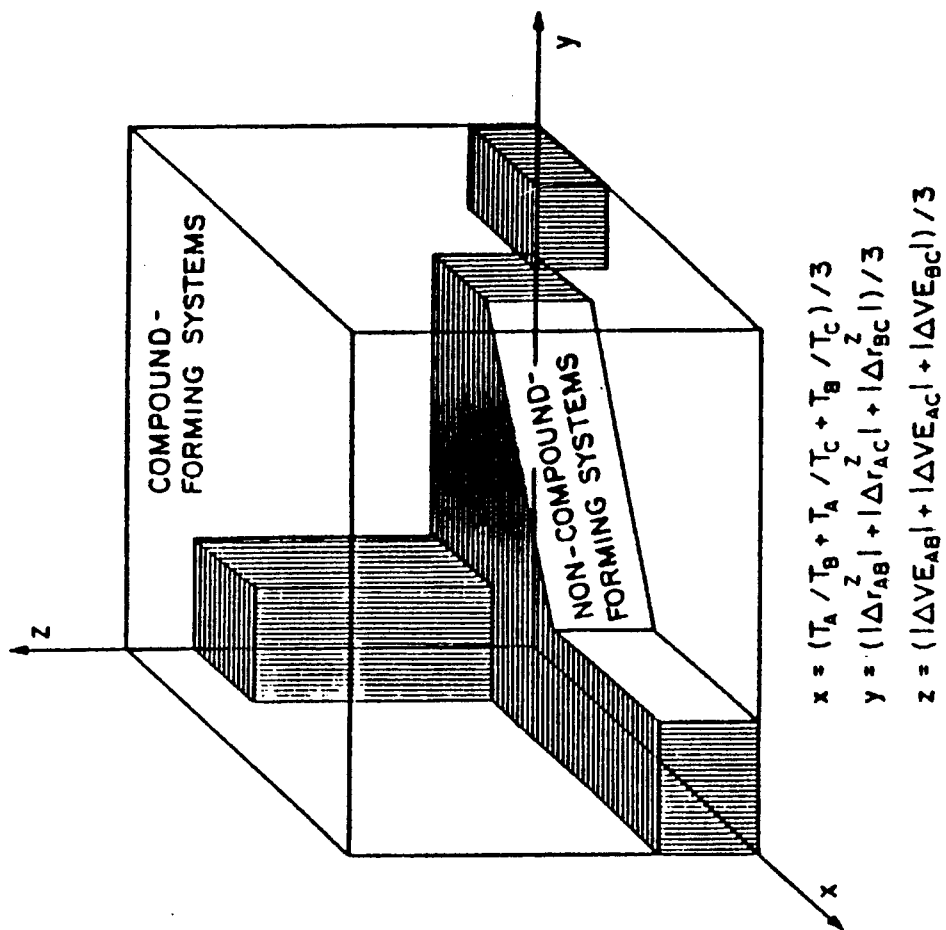
Most frequently occurring AETs

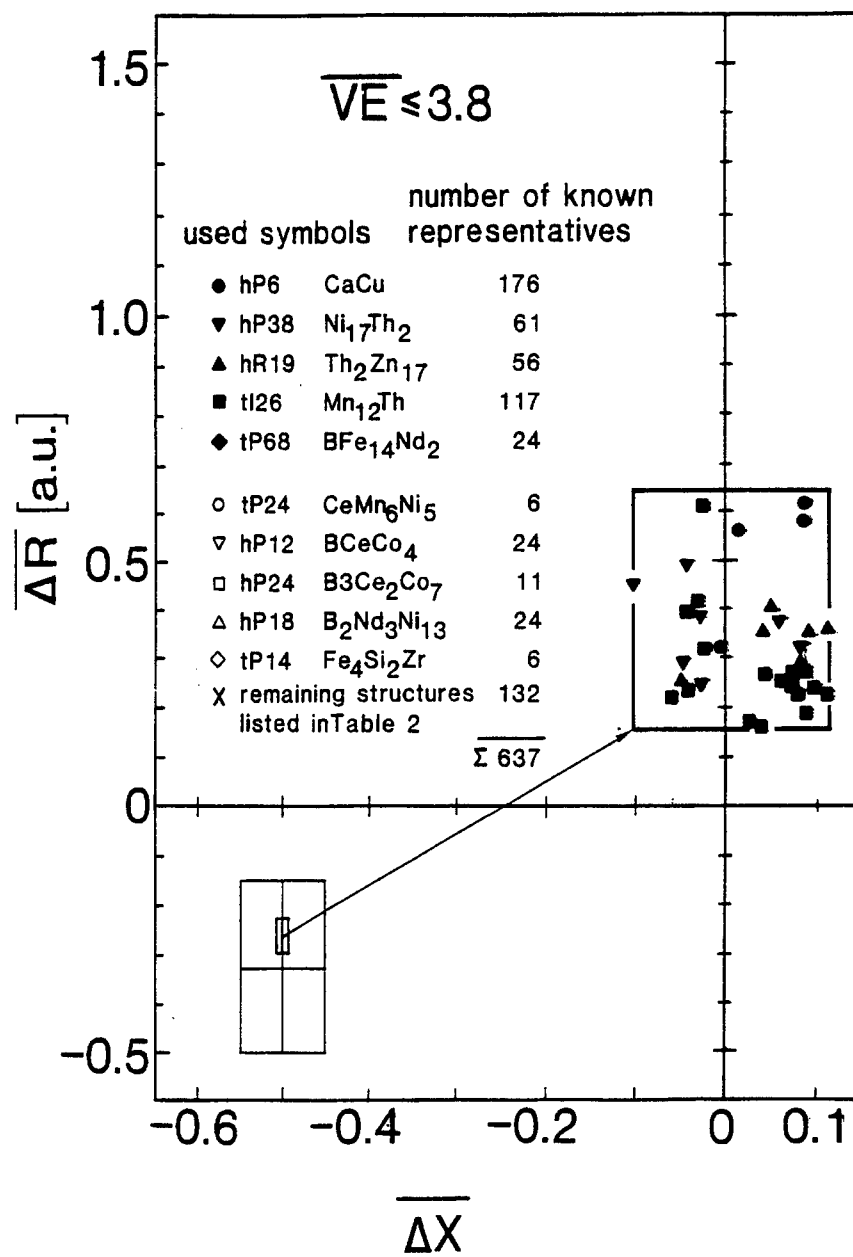


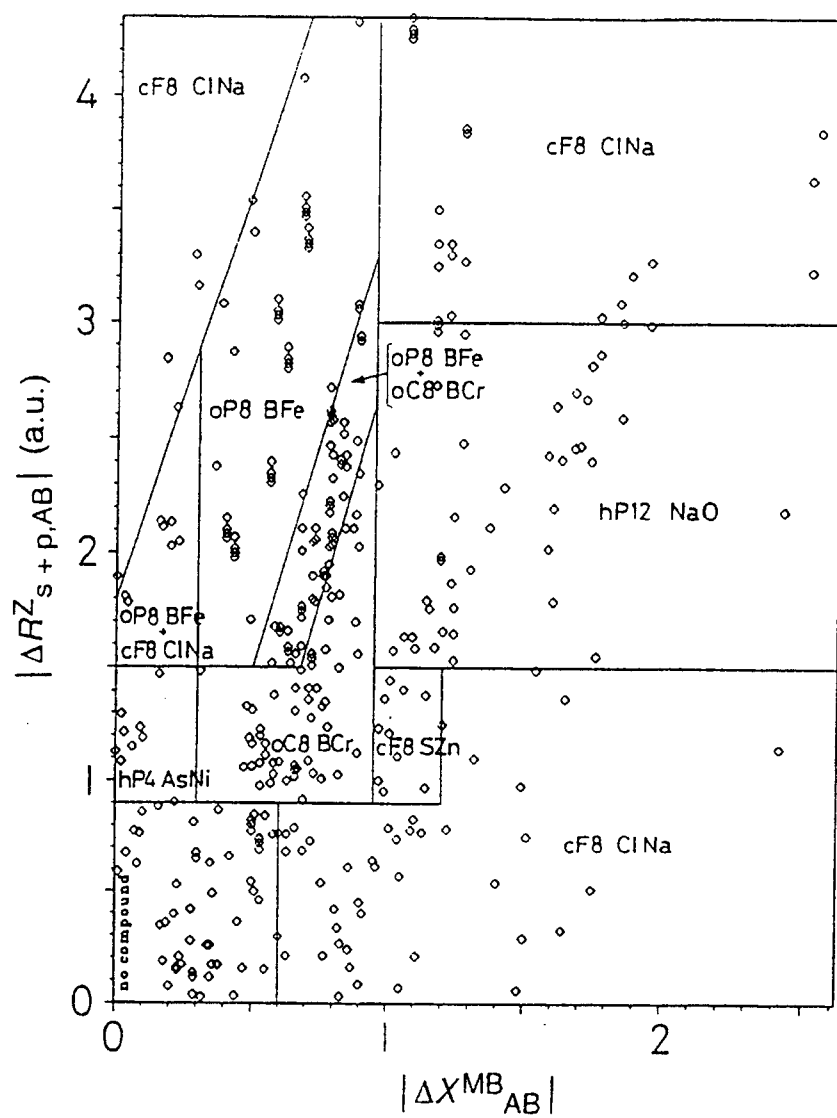
SSD

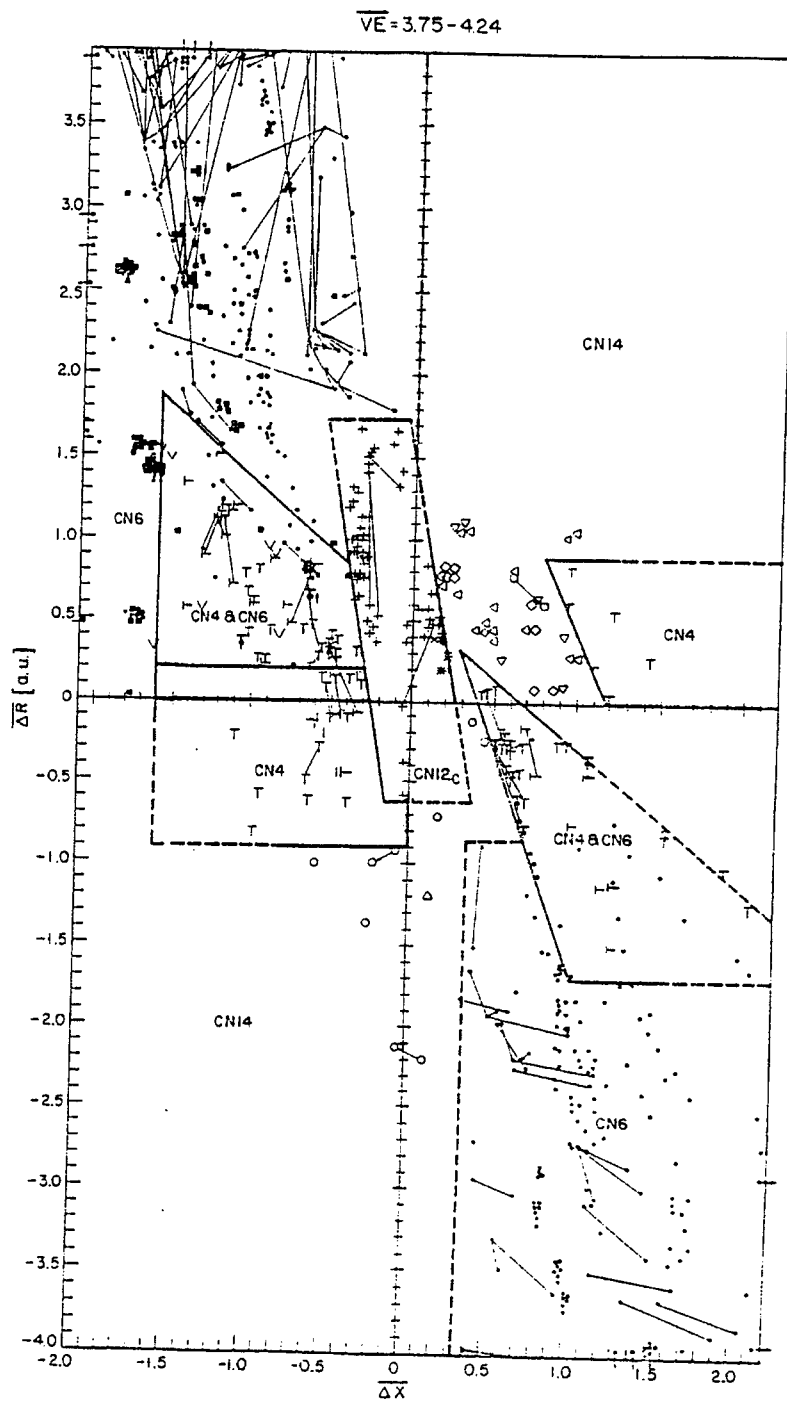


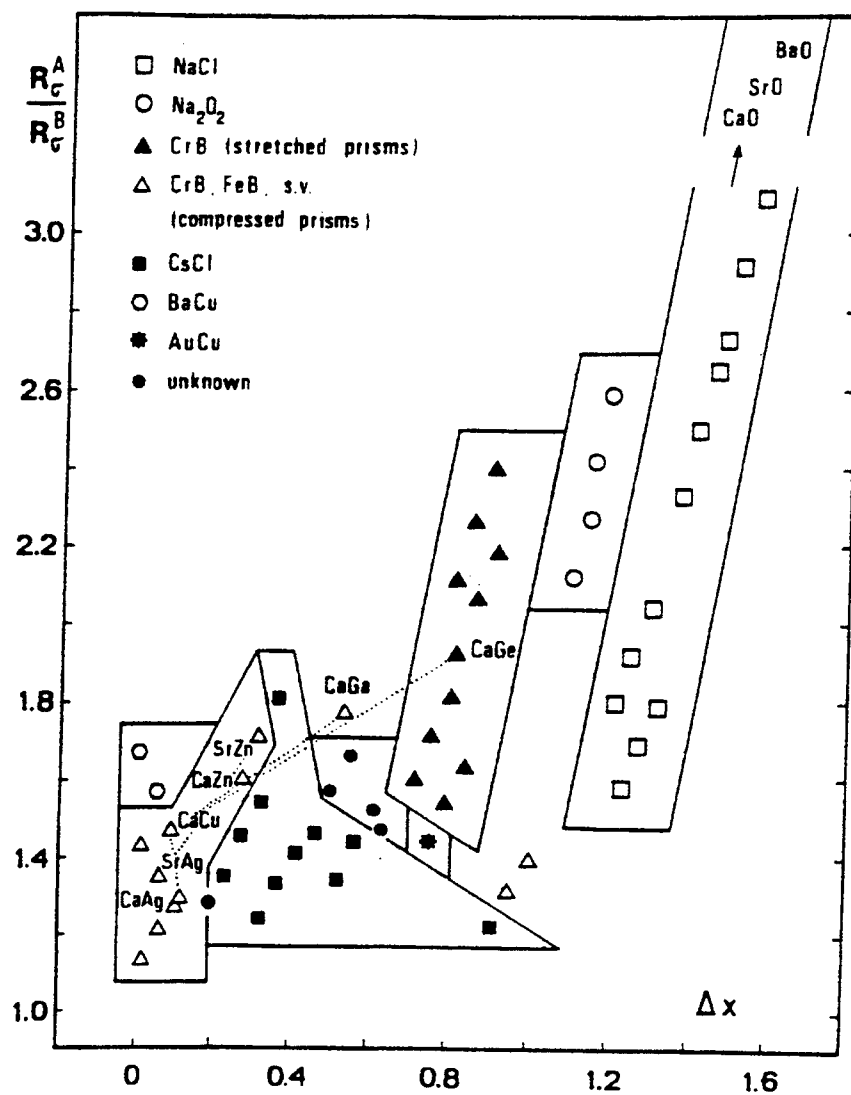


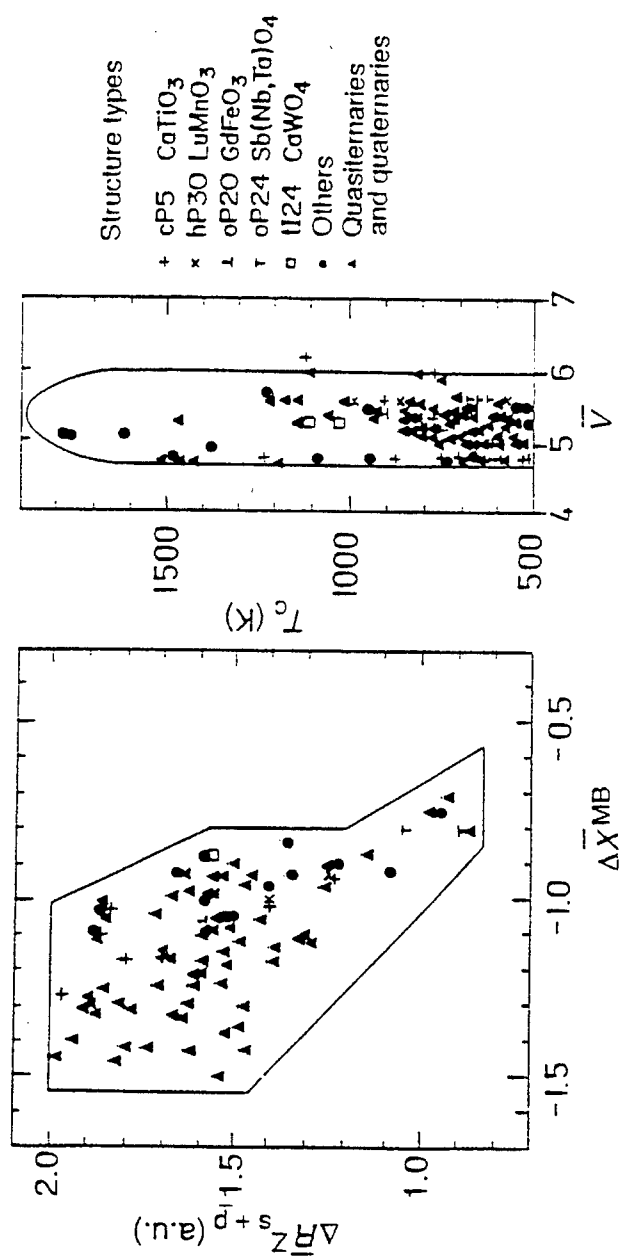


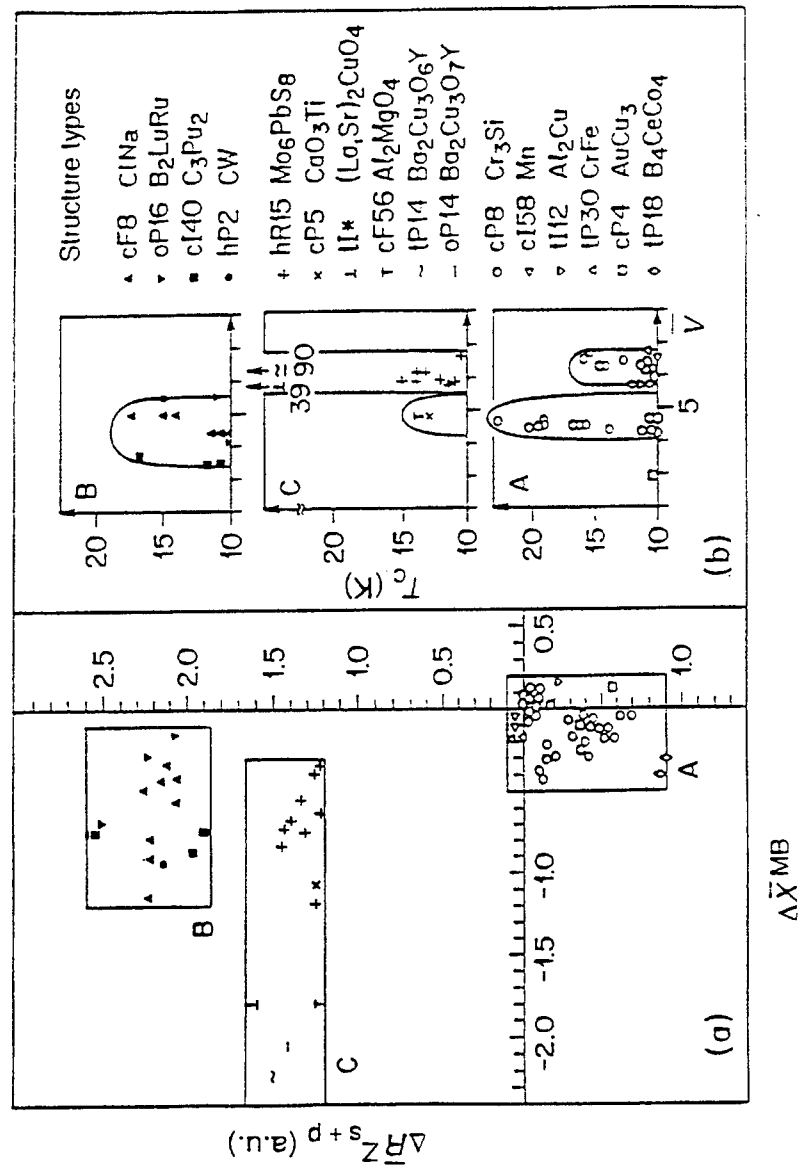












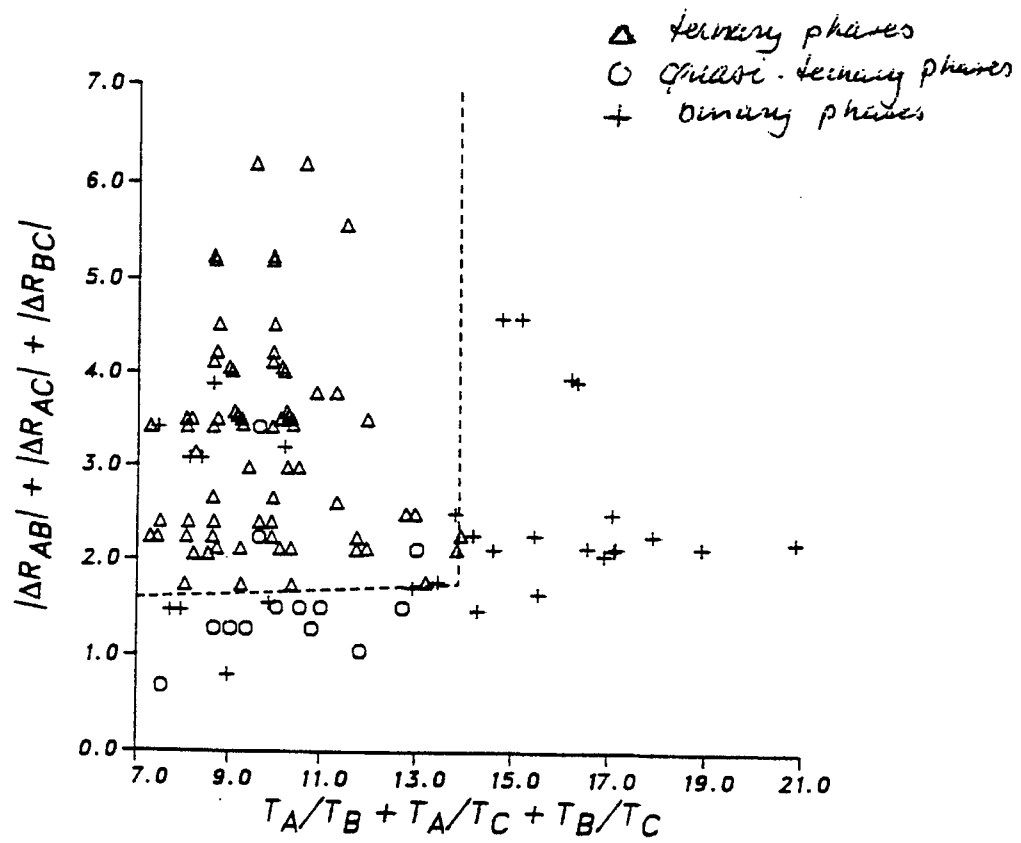
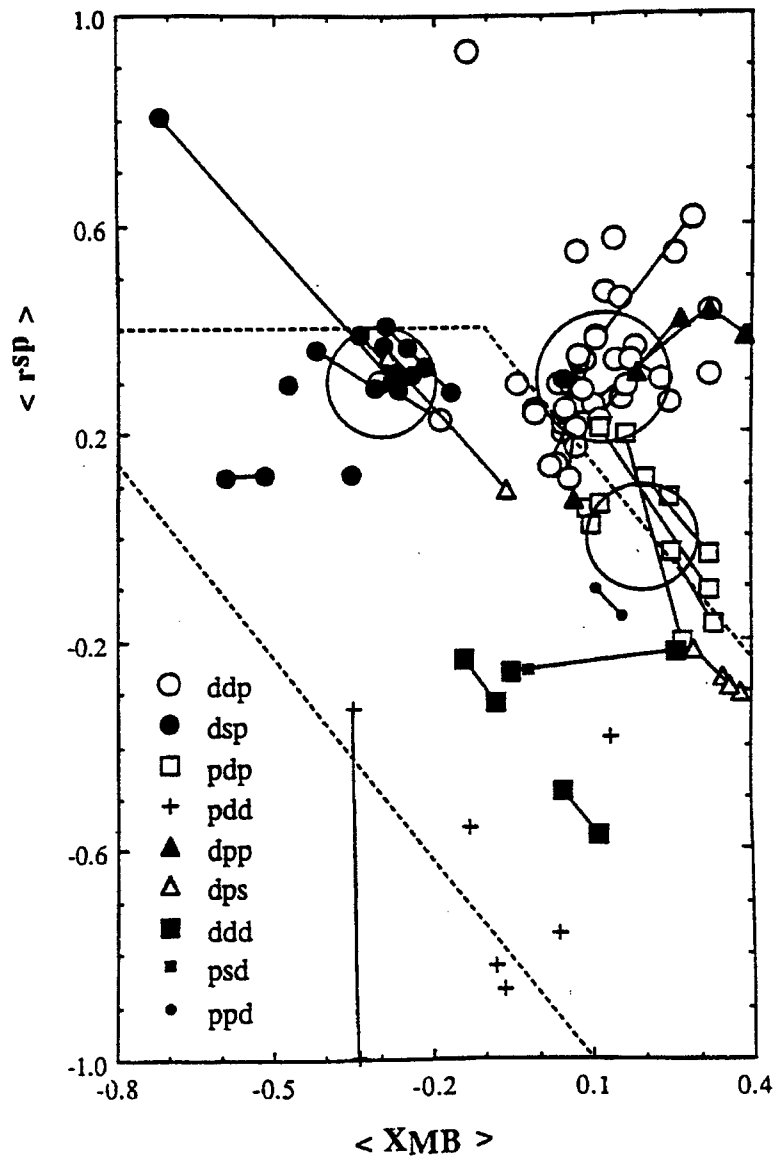
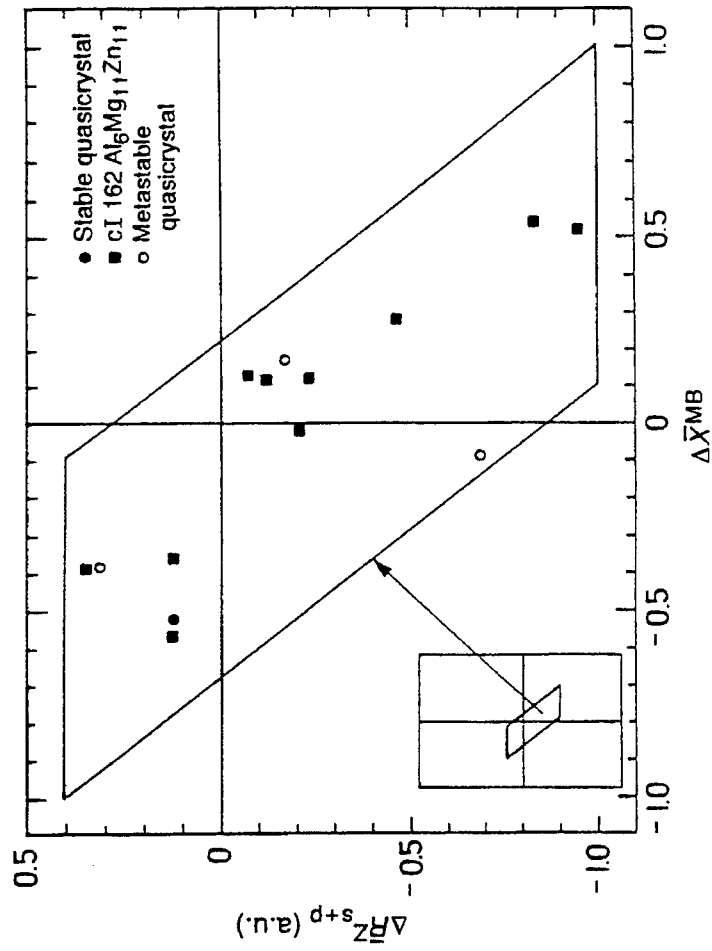


Fig. 5.3

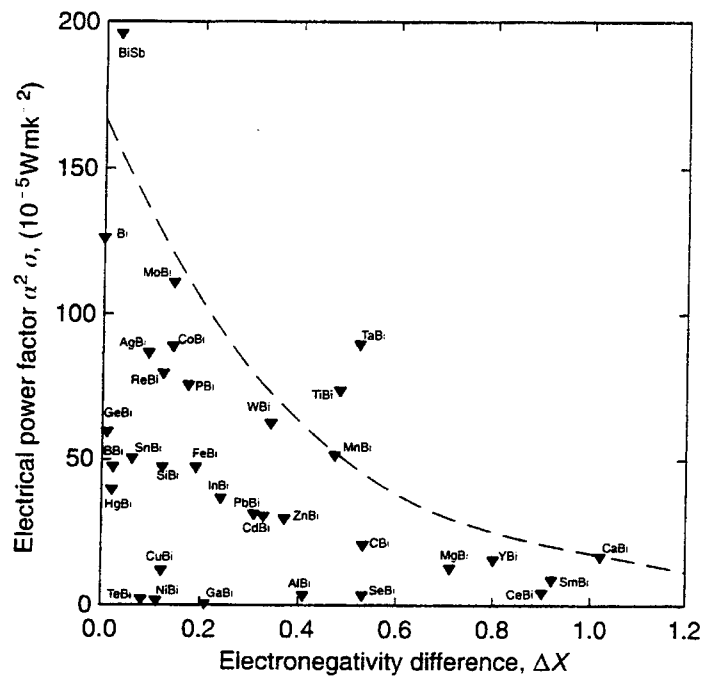


J. Tartar and E.J. Kayastan

J. Mater. Res., Vol. 6, No. 6, Jun 1991

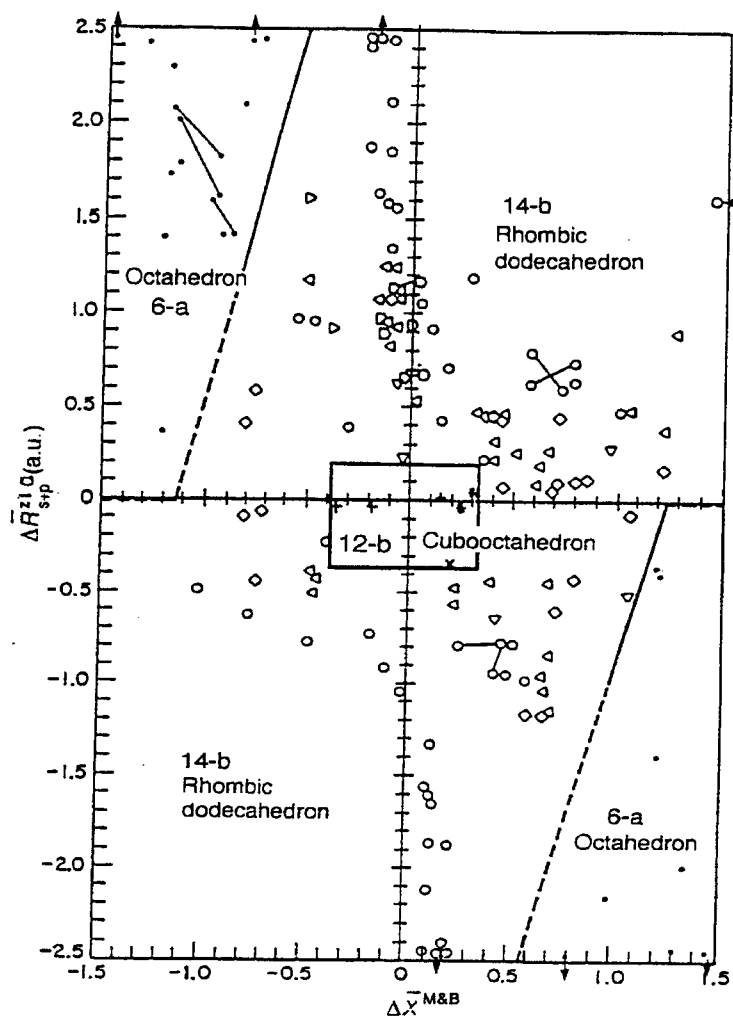


P. Villan

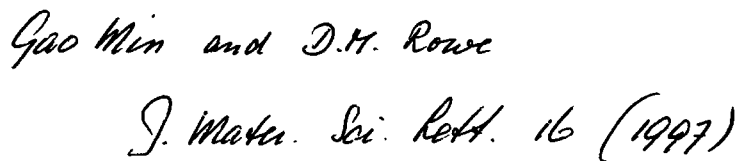


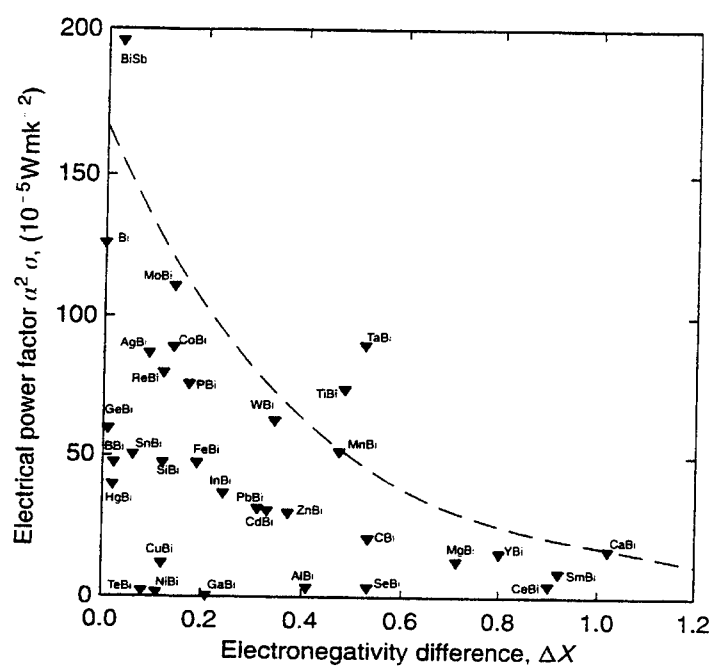
Gao Min and D.H. Rowe

J. Mater. Sci. Lett. 16 (1997)



P. Villars





Geo Min and D.H. Rowe

J. Mater. Sci. Lett. 16 (1997)

Analysis and Visualization of Category Membership Distribution in Multivariate Data

Y.H. Pao*, B.F. Duan*, Y.L. Zhao* and S.R. LeClair**

* Case Western Reserve University
Cleveland, OH 44106

** U.S. Air Force Research Laboratory
WPAFB, OH 45433-7746

This paper reports on some advances in generic data processing procedures with focus on a specific materials discovery and design task. The task is to predict whether a new ternary materials system would be compound forming or not, with the prediction to be based on knowledge of many other known exemplars. The activities and results of three related efforts are described in condensed form in this paper. In one effort, using a combination of clustering and mapping procedures, an accuracy of more than 99 % was attained in predicting the category status [compound forming or not] of new ternary systems. A second effort addressed the question of how to identify redundant or superfluous features. A procedure for identifying the extent of functional dependency amongst features was developed. That procedure can be used to remove redundant features. A third effort addressed the question of how to obtain reduced dimension representations of multivariate data, albeit at the cost of loss of some information. Visualizations of low-dimensional representations can be helpful in building up holistic views of data space, for use in exploration and discovery of new materials

Keywords: Materials systems, compound forming, ternary systems, self-organization, cluster analysis, regional analysis, category membership, membership homogeneity.

Introduction

In the analysis and understanding of large bodies of complex multivariate data, categorization can be of help. Categories serve as large filing folders in the organization of knowledge and data. Category labels serve as filing labels. If the category membership of an object is known, many other attributes of that object can be inferred. The manner in which category membership varies with position in data space can also provide suggestions for discovery of new objects with specific properties.

But it is difficult to understand a large body of high-dimensional data. For example, if a body of data items are presented with ostensibly known category membership labels, it is difficult to portray the distribution of the category membership information throughout the high-dimension data space. Accordingly it is then difficult to judge whether the data information is reasonable, what regions are homogeneous or nearly homogeneous in category membership and what regions are not, and what inferences might be drawn about the likely category membership of other possible new data items. This document describes a suggested systematic procedure for attaining an understanding of large bodies of high-dimensional data. The methodologies include the following:

- (a) Self-organization of data items into clusters on the basis of similarity in high-dimensional data space
- (b) Analysis of clusters in terms of category homogeneity
- (c) A neural-net analysis of the details of category membership distribution within the clusters
- (d) Analysis of clusters in terms of regional distribution of category membership
- (e) Reduced-dimension visualization of distribution of category membership

Item (a) consists of known practice. New insights and practices are introduced in (b) and (c), with item (c) yielding outstanding results, and item (e) is new both in theory and in practice. The work discussed in this paper was motivated by the need for a materials discovery and design task. The objective is to be able to predict whether new compounds might be formed with a set of three atomic elements. Much is known about binary systems but less is known about ternary systems. However, it would be desirable if one were able to allocate any and all ternary systems to one of the two categories, compound-forming or non-compound-forming. Expensive empirical searches for compounds can be avoided if there is guidance from data that it is highly unlikely that the ternary in question would form compounds and new materials, or not.

ANALYSIS AND VISUALIZATION OF CATEGORY MEMBERSHIP DISTRIBUTION IN MULTIVARIATE DATA

Yoh-Han Pao

Baofu Duan

Yanli Zhao

S.R. LeClair

04/21/1999

Case Western Reserve University



Topics Covered

- Prediction of Category Membership of New Materials Systems
- Investigation of Feature Redundancy
- Search for Alternate Features

04/21/1999

Case Western Reserve University



Statement of Task

Predict whether a new ternary materials system will be compound-forming or not, with knowledge of a body of known exemplars.

04/21/1999

Case Western Reserve University



Representation of Ternary Systems

Each element is described by 5 features, i.e., each ternary system is represented by a vector in 15-dimensional space.

Feature Names

- Zunger Radii (Z_r)
- No. of Valence Electrons (VE)
- Melting Temperature (MT)
- Atomic Number (AN)
- Electronegativity (EN)

04/21/1999

Case Western Reserve University



Data Sets Used

<i>Data Set</i>	<i>Population</i>	<i>Forming Status</i>	<i>Used As</i>
Set 1	1067	Non-former	Test
Set 2	2327	Non-former	Training
Set 3	4031	Former	Training
Set 4	244	Mixed	Test
Set 5(*)	4742	Mixed	Training
Set 6(*)	1616	Mixed	Test

* part of (set 2 + set 3)

04/21/1999

Case Western Reserve University



Estimate Category Membership (1)

● Intra-Cluster Neural Net Prediction

- clustering in 15D space, 158 clusters formed
- training one neural net for each cluster
- to predict, first find the nearest cluster, then apply the neural net of that cluster
- advantage: fast, acceptable accuracy
- disadvantage: in case a pattern is on the border of several clusters, the mapping may not be learned properly by only one net.

04/21/1999

Case Western Reserve University



Estimate Category Membership (2)

- **Local K Nearest Neighbor Neural Net Prediction**

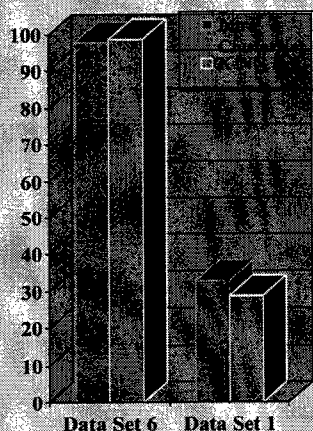
- training neural net based on K nearest patterns
- one specific net for each test pattern
- advantage: mapping is learned for local neighborhood
- disadvantage: time-consuming, computationally expensive

04/21/1999

Case Western Reserve University



Prediction Results



04/21/1999

Case Western Reserve University



- Trained on Set 5
- Accuracy of 98.9% is obtained for Set 6
- Only about 30% accurate for Set 1
- High accuracy for Set 4

Nonlinear Functional Correlation Amongst Features

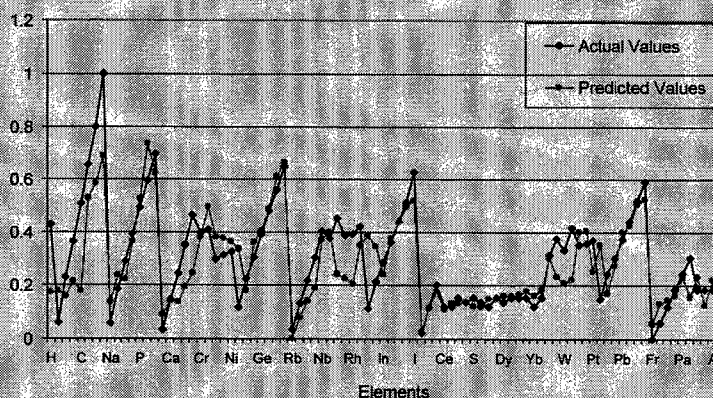
- If a feature is functionally correlated with other features, its value could be predicted
- Learn a mapping from $N-1$ features to the N th feature
- Training on 1/3 of the elements, 1/3 for generalization, and 1/3 for test.

04/21/1999

Case Western Reserve University



Predict Electronegativity

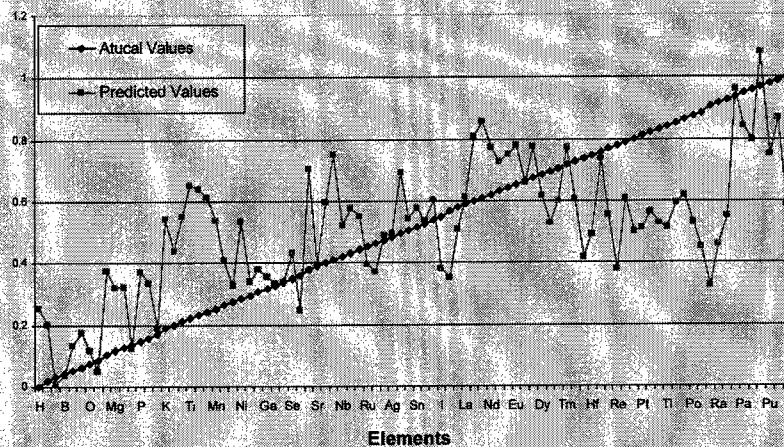


04/21/1999

Case Western Reserve University



Predict Atomic Number

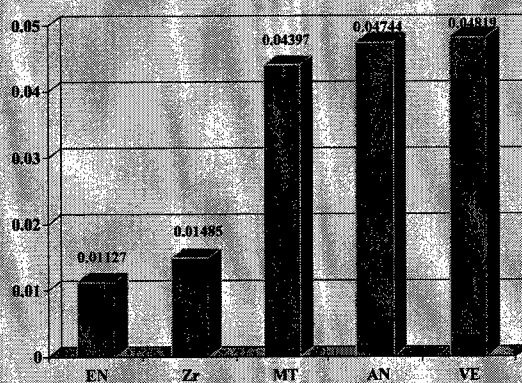


04/21/1999

Case Western Reserve University



Variance of Error in Predicting Features



$$Var = \frac{1}{90} \sum_{i=1}^{90} (f_i - \tilde{f}_i)^2 \quad f_i : \text{actual value}, \tilde{f}_i : \text{predicted value}$$

04/21/1999

Case Western Reserve University



Question:

Is the functionally correlated feature
superfluous in terms of predicting category
membership?

04/21/1999

Case Western Reserve University



Identification of Superfluous Features

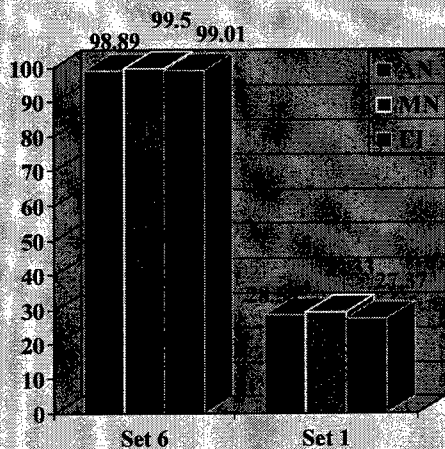
Features Excluded	None	AN+EN	AN+Z _i	AN+VE	AN+MT
Error	18	21	23	30	34
Correct %	98.89	98.7	98.58	98.14	97.79
Features Excluded	AN	MT	VE	Z _i	EN
Error	18	23	23	24	27
Correct %	98.89	98.58	98.58	98.51	98.33

04/21/1999

Case Western Reserve University



Search for Alternate Features



04/21/1999

- With no specific preference, Mendeleev Number(MN), Energy Ionization 1st (EI) were used as alternate feature for Atomic Number (AN)
- Knowledge is needed to guide the search for more alternate features

Case Western Reserve University



Comments on Results

- Local net for predicting category membership is an efficient and effective approach;
- Functional correlation between features may not be useless in predicting category membership for ternary systems ;
- Prediction of category membership might be possible with less than five features per element;
- Knowledge of materials systems is needed to guide the search for alternate features.

04/21/1999

Case Western Reserve University



Whitney Reduction Networks for Process Discovery

Mark E. Oxley

Department of Mathematics and Statistics
Air Force Institute of Technology
2950 P Street, Wright-Patterson AFB, OH 45433-7765
Email: moxley@afit.af.mil

ABSTRACT

One of the best known methods for dimensionality reduction is the Karhunen-Loeve transform, or principal component analysis (also known as proper orthogonal decomposition or singular value decomposition). This empirical linear transformation is appealing given it is based on solving an eigen-vector problem which produces an optimal spanning subspace. Alternatively, empirical nonlinear transformations may be generated, for example, by computational neural networks or radial basis function expansions. These adaptive mappings may be determined by solving both smooth and combinatorial optimization problems. Such nonlinear mappings have the advantage that they may be used to efficiently model manifolds as opposed to subspaces.

This talk will present an application of Whitney's embedding theorem to the data reduction problem and will introduce a new reduction technique motivated, in part, by a constructive proof of the theorem. In this setting, we introduce the notion of a "good projection". We show it is useful to optimize empirical projections with respect to their inverses, i.e., these should be well-conditioned. One possibility is computation of the singular vectors of the secants of the data. This may be improved upon by using an adaptive algorithm. A method for constructing the nonlinear inverse of the projection and a discussion of its properties will also be presented. Finally, well-known methods of data reduction are compared with our approach within the context of Whitney's Theorem.

Introduction

Complicated physical processes which have dissipative dynamical evolution

- . exhibit self-organizing behavior
- . are naturally described by low-dimensional models
- . have data that tends to cluster in small volumes of the total space

It is within these subspaces, or submanifolds, that low-dimensional description may be obtained.

Dimensionality Reduction

Let $a \in A \subset R^n$ and G be a dimensionality reducing mapping, i.e.,

$$G : A \rightarrow B \subset R^m$$

$$b = G(a)$$

Let H be a reconstruction mapping

$$H : B \rightarrow A$$

$$a = H(b)$$

Empirically determine approximations \tilde{G}, \tilde{H} to G, H such that

$$\left\langle \left\| a - \tilde{H} \circ \tilde{G}(a) \right\|^2 \right\rangle = \min$$

See Kirby (1999) for extensive discussion.

Whitney Reduction

Easy Whitney Embedding Theorem

Let M be a compact m -dimensional submanifold of R^q where $q > 2m+1$.

Then there is an embedding of M in R^{2m+1} .

This theorem says:

- Linear reduction permissible (almost any)
- Reduction dimension attained $d = 2m + 1$
- Reconstruction nonlinear, but domain given.

Additionally, we seek to construct mappings of *minimum complexity*.

Process Analysis

Let f denote the process (known) with input set I , control parameters Ω , and output set O .

$$f : I \times \Omega \rightarrow O$$

Given a reduction mapping $G : O \rightarrow \tilde{O}$ and reduced model \tilde{f} now

$$\tilde{f} = G \circ f : I \times \Omega \rightarrow \tilde{O}$$

and

$$f = H \circ \tilde{f}$$

Choice of Mappings

- H and G linear
- H and G nonlinear
- H linear, G nonlinear
- H nonlinear, G linear

Globally or locally

H and/or G analytical, empirical or mixed

Use... ..

- Neural networks
- Radial basis functions
- PCA (i.e., SVD)
- Other

Example: Bottleneck Neural Networks

Reduction mapping G is nonlinear

Reconstruction mapping H is nonlinear.

At least 3 hidden layers in feedforward MLP; nonlinear optimization problem.

(See Kirby and Miranda, 1999)

Process Discovery

Let f denote an unknown process

Given data in I , Ω , O determine an approximation to f .

Let $A = I \times \Omega \times O$ choose reduction mapping

$$G : A \rightarrow B$$

thus we work in a *smaller* space B .

Conclusions

1. Whitney's theorem guarantees that we can work in a smaller space.
2. Working in this smaller space will be less complex to discover the unknown process.

References

David Broomhead and M. Kirby (1998), *New Approach for Dimensionality Reduction: Theory and Algorithms*, to appear SIAM J. of Applied Mathematics.

David Broomhead and M. Kirby (1998),
The Whitney Reduction Network: a method for computing autoassociative graphs,
 (submitted to Neural Computation)

Kirby (1999) *Dimensionality Reduction*, to appear, Wiley & Sons.

M. Kirby and R. Miranda (1999), *Empirical Dynamical System Reduction } : Global Nonlinear Transformations*,
 {In: Semi-Analytic Methods for the Navier-Stokes Equations (Montreal, 1995)},
 {K. Coughlin}, Vol 20, pp41-64,
 CRM Proc. Lecture Notes, Amer. Math. Soc., Providence, RI

M. Kirby, L. Sirovich, (1990),
Application of the Karhunen-Loeve procedure for the characterization of human faces, IEEE Trans. PAMI, Vol 12, No. 1, 103.

Intelligent Materials Processing by Hyperspace Data Mining

Nianyi Chen, Dongping Daniel Zhu and Wenhua Wang*

Zaptron Systems, Inc., Mountain View, CA USA

Email: dan@zaptron.com

*Salomon Smith Barney, 250 West Street 8th Floor, New York, NY 10013 USA

ABSTRACT

This paper discusses application of hyperspace data mining to materials manufacturing. We introduce an innovative hyperspace method whereby data are separated into subspaces, features selected according to data patterns, and control rendered in the original feature space. This technique has three major advantages: no equipment is added, no experiment is needed, and no interruption occurs to production. A number of proprietary algorithms have been built into a software product, MasterMiner™, for use in materials design and manufacturing. Examples are given to show the efficacy of the proposed method and MasterMiner tool.

INTRODUCTION

In spite of a 5,000 year long history of developing, manufacturing and using materials, drugs and chemical products, we continue to search for new materials and substances to meet the ever-increasing needs in industry, agriculture, health, and national defense. Historically, the exploration and development of new materials (substances) has been based on experimental or trial-and-error methods. The more recent trend in materials and molecular design has focused on studies of the relationships between material structure and material property by quantum chemistry. Today, more sophisticated methods, such as artificial intelligence, data bases and computer information processing, are being used to build expert systems to narrow the search for desired information. The computer-aided methods have been the main stream in the field of advanced materials design. As the latest method in materials research and process control, intelligent processing of materials is adopted in manufacturing new materials.

Today, materials design and manufacturing is not only a hot research topic at research institutes, but also an important business in industry. A completely new methodology is being developed in the research and development of advanced materials and substances. In general, there are three levels of computerized materials design: Level 1. Use quantum chemistry, solid physics, and structural chemistry to study relationships between material microstructure and property; Level 2. Explore the development of new alloys, ceramics and semiconductor materials by using phase diagrams, thermodynamics and kinetics; Level 3. Apply pattern recognition, data mining, neural nets, and genetic algorithms to optimize the manufacturing, processing and property of materials. Very often this is accomplished with the aid of data base, knowledge base and knowledge discovery. The method reported herein falls in Level 3.

Data mining [1] [5] [6] is the process of discovering meaningful new correlations, patterns and trends by sifting through large amounts of data stored in repositories, using pattern recognition technologies as well as statistical and mathematical techniques. Data mining in fact is an optimization technique that has found practical applications in many industries, including materials design and manufacturing, drug screening and production, steel making, power generators, petro-chemical, and operations management [1] [2].

The most important part in material data mining is the capability to separate data samples in a multi-dimensional hyperspace spanned by a number features. The task of data mining is to determine the data pattern from a given set of data. Once a data pattern has been recognized, an artificial neural network can be trained to find the mathematical model for the pattern. Then a genetic algorithm can be developed to search for the global optimum. Finally, the knowledge obtained by the expert system, including the mapping graphics by data mining, criterion equations, trained neural nets and the optimum boundary, is stored in a knowledge base for later use.

PROBLEM BACKGROUND

By nature, a material design problem is an optimization problem, and methods in pattern recognition and data mining can be used to offer effective solutions. Most pattern recognition methods are based on the computerized recognition of the multidimensional graphs (or their two-dimensional projections) of the distribution of samples from different classes in a multidimensional space. Independent variables (often called system input, features or factors) influencing the target (dependent variable or system output) are used to span a multidimensional space.

We can describe samples of different classes as points with different symbols in these spaces. Various pattern recognition methods can be used to "recognize" the patterns shown in the graph of *distribution zones* of different samples. In this way, a model, qualitative or quantitative, can be obtained that describes the relationship (or regularity) among targets and factors. If we adjust criterion of classification, semi-quantitative models describing the regularities can be found, if noise is not too strong. Unlike regression methods (linear regression, nonlinear regression, logistic regression, etc.) or the artificial neural networks (ANN) [4] that provide *quantitative* solutions, pattern recognition methods often provide *semi-quantitative* or *qualitative* solutions to classification. This is of course a limitation of pattern recognition methods. However, this is not always a disadvantage, because many data sets exhibit strong noise, and a quantitative calculation would be *too precise* to present them. Besides, practical problems in many cases are of the "yes or no" type, and pattern recognition is especially suited to offering adequate solutions to them. For example, a problem may be "whether the fault will occur or not", or "whether an intermetallic compound will form or not."

A number of common pattern recognition methods have been built into MasterMiner software. They include Principal Component Analysis method (PCA), Fisher method (Fisher), Partial Least Square method (PLS) [3], Sphere-Linear Mapping method (LMAP), Envelope method (Envelope), Map Recognition method (MREC), and Box-Enclosing method (BOX). PCA, PLS and Fisher methods are traditional pattern recognition methods, and their principles are described in standard textbooks. In general, limited separation is achieved by traditional PCA or Fisher method when the data exhibit strong non-linearity. The other four methods listed above are developed specially for processing complicated data sets. In the sphere-linear mapping (LMAP) method, computation starts by moving the origin of the initial data coordinate system to the center of sample points of class "1", followed by finding a hyper-ellipsoid to enclose all sample points of class "1". By a *whiten transformation*, this hyper-ellipsoid is changed to a hyper-sphere. The multidimensional graph of the data points is then projected onto a series of two-dimensional planes to form a number of 2-dimensional maps on the computer screen.

A HYPERSPACE DATA MINING METHOD

Data Separability

The data separability test of MasterMiner is designed to investigate the possibility of separating data points from different populations or clusters in a hyper-space. If the data are separable, it may be possible to build a mathematical model to describe the system under study. Otherwise, a good model can not be built from the give data and more data or data processing is needed. MREC (map recognition by hidden projection) is a method used to choose the "best" projection map from a series of hidden projection maps. Here "best" implies that the separation of sample points of different classes is better than those obtained from other maps. In MasterMiner, MREC is used together with the "auto-square" method to provide the so-called MREC-auto-square solution. In each projection map, sample points of class "1" are automatically enclosed by a square frame (as shown in Fig. 1), and a new data set is formed that contains only sample points within this "auto-square." This new data set will be used to in model building for diagnosis or failure recognition.

After this auto-box operation, a second MREC is performed on the new data set to obtain a new "best" projection map where sample points remaining in the auto-square are separated better into 2 classes. After series of such hidden projections, a complete (close to 100%) separation of sample points into different classes could be realized. It has been shown that MREC-auto-square method is much more powerful than the traditional patter recognition method.

The physical meaning of MREC-auto-square method is explained as follows: each "auto-square" represents actually a square "tunnel" in the original multidimensional space, and several such tunnels would form a *hyper-polyhedron* in this space by intersection of all tunnels, as shown in Fig. 1. This hyper-polyhedron, enclosing all "1" sample points, defines an optimal zone in the multidimensional space, if all or most of the samples of class "2" are separated from this zone because they are located outside of this hyper-polyhedron.

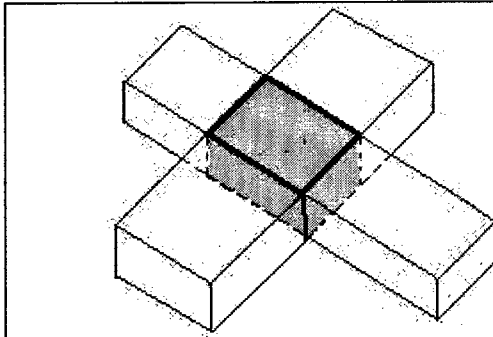


Fig.1 Polyhedron formed by square tunnels.

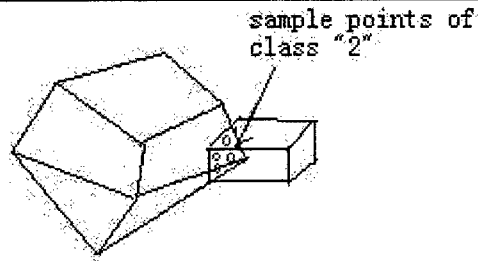


Fig.2 Principle of the Envelope-Box Method

Back Mapping

Since the hidden projection by MREC transforms data from the original measurement (or feature) space into a number of other orthogonal spaces, we need to back map the transformed data into the original feature space to derive mathematical models for practical use. A method called PCB (principal component backing) [2] has been developed whereby a point representing an unknown sample from a low-dimensional principal component subspace is back-projected to the high-dimensional space of original features. In PCB, non-linear inverse mapping and linear inverse mapping are used to obtain an accurate solution to predicting sample points in the optimal region.

Let \mathbf{X} be a standard training set with n samples and m features, and let \mathbf{Y} be the sample set in the PC (principal component) space corresponding to \mathbf{X} in the original feature space. We have $\mathbf{Y} = \mathbf{XC}$, where the columns of \mathbf{C} are the eigenvectors of the covariance matrix \mathbf{D} ($\mathbf{D} = \mathbf{X}^T \mathbf{X}$) for the training set \mathbf{X} . The 2-d subspace of PCs consisting of \mathbf{C}_u and \mathbf{C}_v is defined as the main map where samples are assumed to be completely classified. A point P in the main map represents the unknown sample, and it is described by two variables, y_{pu} and y_{pv} , respectively. In general, p is expected to be an optimal sample if its neighbors are optimal points. In order to back transform an unknown sample point to the original space, i.e., to find \mathbf{X}_p^* , one has to determine its boundary conditions; otherwise, an uncertain solution will occur. Two types of conditions are proposed for PCB: non-linear inverse mapping and linear inverse mapping.

In non-linear inverse mapping (NLIM), let the error function E be defined as

$$E = \left(\sum_{j=1}^n d_{pj} \right)^{-1} \sum_{j=1}^n \frac{(d_{pj} - d_{pj}^*)^2}{d_{pj}}$$

where

$$d_{pj}^* = \left(\sum_{k=1}^m (x_{pk} - x_{jk})^2 \right)^{1/2}$$

$$d_{pj} = [(y_{pu} - y_{ju})^2 - (y_{pv} - y_{jv})^2]^{1/2}$$

Here d_{pj} is the distance from the unknown sample represented by p to all known samples in the subspace defined by PCs two coordinates, u and v , and is the same distance in the original feature space. Non-linear optimization method is utilized to compute the values of z_{pk} that minimizes the error function E . The solution using NLIM boundary condition is only an approximation, because the parameters obtained in this way depend to some extent on the trial coordinates in the original space.

Linear Inverse Mapping (LIM): Besides the 2-d subspace of PCs consisting of C_u and C_v , there exists an $(m-2)$ -dimensional subspace of PCs consisting of C_i ($i = 1, 2, \dots, m$ and $i \neq u, v$), since C is derived from the covariance matrix D ($m \times m$). When the projection of point p , which is described by y_{pv} and y_{pu} in the main map, are determined with y_{pi} ($i = 1, 2, \dots, m$ and $i \neq u, v$) in the $(m-2)$ -dimensional subspace, a set of simultaneous *linear* equations can be obtained as

$$y_{pk} = \sum_{j=1}^n C_{jk} x_{pj}$$

Where $k = 1, 2, \dots, m$, and the set of linear equations can be solved for the parameters of the unknown sample point corresponding to point p . The linear inverse mapping will always produce an exact solution. As to the projection of point p , in general one can let point p be at the center of the region containing the largest number of known optimal samples so that the unknown sample has properties similar to its neighbors, the known optimal samples.

Auto-box for Concave Polyhedron

Since both the MREC and envelope method can only form a *convex hyper-polyhedron*, they can not separate the sample points of different classes if the distribution zone is not a *convex* but a *concave* polyhedron. In these cases, another technique, the BOX method shown in Fig. 2, can be used as an effective solution. When there are some sample points of class "2" still remaining within the hyper-polyhedron after applying the envelope method, a smallest multidimensional "box" will be built by MasterMiner to enclose all sample points of class "2" (or class "1").

We can build a box within the hyper-polyhedron in an effort to enclose all "2" sample points. If the box so built contains no sample point of class "1" (or only very few points of class "1"), rather good separation will be achieved by throwing away those sample points (mostly of type "2") in the box. In fact, BOX of MasterMiner is a tool for virtual data mining. It can also be used to separate sample points of class "2" from that of class "1". BOX method itself has limited use, but it becomes a very useful tool in virtually mining the data space when combined with the MREC or Envelope geometrical method.

Factor Selection - not a text-book approach

Analysis is performed in the m -dimensional factor (or feature) space. Any system model must be built on the selected factors that represent the operating rules of the system. Therefore, selection of an appropriate set of factors is very important in data mining. Inappropriate selection of factors may lead to unnecessary complexity and introduce noise to the data. There are two general approaches to feature selection. One approach is to use the *first principle* method to study the physical, chemical or electrical properties and perform experiment and take measurement to diagnose major factors. The second approach is to use pattern recognition and optimization techniques to find the relationships among various factors from history data.

Strictly controlled factors - many important factors are already under close control and vary very little in production process, and they should not be considered as important factors. Pay more attention to other factors.

Common vs. specific factors - common factors from text book or common sense knowledge may not apply to a specific process in a specific case. Attention should be directed to finding those factors that are specific to the process under study. These specific factors are more important than the common factors.

Evolutionary factors - even with the same process in the same plant, the priority of the identified factors may change. While one factor may be the deciding factor in solving the bottleneck problem, it may be less important than some other factors in solving the product quality problem of the same plant.

Minimum set of effective factors: this is the min set of factors that can be used to represent the system under study. An interactive and iterative technique has been developed in MasterMiner to identify this factor set, which has been proved highly effective and efficient in many industrial applications.

Factor Multiplicity - a concept borrowed from molecular chemistry that describes the multiplicity in the phase change of a substance. In an optimization problem, $Y = f(X_1, X_2, X_3, \dots, X_i, \dots, X_n)$, factors $X_1, X_2, X_3, \dots, X_i, \dots, X_n$ and target Y are interchangeable. For instance, $X_i = g(X_1, X_2, X_3, \dots, X_i, \dots, X_n)$ describes another optimization problem for the same system. This means that factors are not fixed in a problem, and the challenge is to identify the best factors for one specific problem using an efficient and effective method in addition to expert's knowledge and experience.

EXAMPLES OF DATA MINING IN MATERIAL MANUFACTURING

Case 1 Optimization of Synthetic Rubber Production -- Butadiene Rubber

Background: We used MasterMiner to process data from a rubber factory to improve product quality. The performance of butadiene rubber is, in general, determined by the *ML* value of the product of polymerization. *ML* is a parameter related to the molecular weight distribution of polymeric products. The production requires *ML* value to be in the range of 43 - 47. A rubber with too high *ML* value has low elasticity, and one with too low *ML* value has low tensile strength. Since polymerization involves chemical reactions, heat transfer, mass transfer and fluid flow, it is impossible to find an effective mathematical model from "first principles". The objective is to use historical records to build an empirical relationship, since exact modeling is impossible. The *ML* values from historical data were divided into 3 classes: class 1 for $ML < 43$, class 2 for ML within the interval $[43, 47]$, and class 3 for $ML > 47$. A total of 45 factors were identified that could possibly affect the molecular weight *ML* of the material. These include temperature, flow rate, feed of catalyst, feed of solvent oil, etc. After data mining, we discovered 5 factors $\{Z_1, Z_2, Z_3, Z_4, Z_5\}$, called *principle factors*, that have significant effect, as shown in the following table:

Factor	Property
Z1	feed of butadiene
Z2	feed of solvent oil
Z3	feed of catalyst
Z4	temperature of feed
Z5	temperature at lower part of first reactor

Results: Using MasterMiner, we successfully built a reliable mathematical model for the *ML* value in the 5-dimensional hyperspace spanned by $\{Z_1, Z_2, Z_3, Z_4, Z_5\}$, as follows

$$ML = 2.111 - 0.661Z_1 - 0.00636Z_2 + 0.03737Z_3 + 0.01255Z_4 - 0.02397Z_5$$

This model was used to control the material manufacturing with good results. The rate of on-spec butadiene rubber (of good quality) was increased from 71% to 95.2%, and the production yield rose from 89% to 93%. Net annual revenue in the factory was increased by US \$0.25 million.

Case 2 Predictive Control of Chaotic Processes

Often a process appears to be unreproducible due to *chaotic processing*, i.e., in producing ultra fine Al-powder for PTC, the particle size is inconsistent from batch to batch since reproducibility is not very good. By data mining using MasterMiner software, it was found that the particle size distribution is related to the UV spectra change of the aqueous solution before precipitation of aluminum hydroxide. Based on this relationship, we developed an effective method to improve the quality of the ultra-fine Al-powder. By monitoring the UV spectra before precipitation, we can predict if the particle size of a batch is "too fine" or "too coarse," and then adjust the pH value of the solution to within specification. In this way, the particle size distribution becomes more uniform. In one factory, inconsistent particle size was observed in PTC ceramic production with an on-spec product rate below 60%. A method was needed to control product quality. The material used is a ultra-fine Al_2O_3 powder. The chemical reaction is



Control Process: add acid or base to control (speed up or slower) the above induction process, or change the cooling rate to alter the manufacturing process. Heated $\text{Al}(\text{OH})_3$ powder was formed with a particle size distribution near Gaussian. The Al_2O_3 powder is then formed and used in the control system. A method was designed to discover a relationship wherein a violet light at wavelength 2800 Å is applied to the material being formed to measure transparency variations. MasterMiner provides a predictive control solution to this manufacturing process by predicting product quality 30 min before finish from measurements of the resistance curve of a Al_2O_3 blob being formed. Control is achieved by changing cooling rate to control the final resistance at 60 minutes. Using this method, product quality improved from 60 to 100% in 500 tests.

Case 3 Energy Saving in Aluminum Production

Background: In aluminum production at one aluminum factory, 5,000 electrolytic cells are used in an aluminum leaching process. An electric current of 100,000A is applied to each of the cells. The anode of the cells is made of an electro-conducting alloy (Fe-P-Si-Mn-S). High electric power consumption was caused by the high resistance of micro-cracks formed in the cell alloy during production. The voltage drop across each cell is 0.4-0.5 volts, and the total power consumption by the 5,000 cells in the aluminum production line is about $0.5 \times 100,000 \times 5,000 = 250,000,000$ watts = 250,000 kW. Therefore, the goal is to reduce power consumption by reducing the micro-cracks in the alloy. History data from operations were used to build a feasible model that can be used to control and reduce the formation of micro-cracks in the cell alloy. 5 major factors have identified that control the formation of the micro-cracks on these cells. MasterMiner offered the following advisory to reduce the formation of the micro-cracks: Mn (a1) increased, Si (a2) reduced, C (a3) no change, S (a4) no significant influence, and P (a5) no change. At the end, the cracks were significantly reduced, and the factory saved electricity by 3,000,000 kWh per year

CONCLUSION

The proposed hyperspace data mining is very effective for modeling complex chemical systems, involving heat transfer, mass transfer, fluid flow and chemical reactions. Based on the built models, optimization can be realized by intelligent control. It can optimize many materials preparation processes, including high temperature (Tc) superconductors, ceramic semiconductors (Ga, In), (P, As, Sb) film by MOCVD method, rare earth-containing phosphor, alloy steels, etc. The method is also applicable to optimize etching, VCD and other IC processes in semiconductor production. Industries where MasterMiner™ has been applied are:

- Chemical industry: quality improvement of synthetic rubber, fiber, plastics and fine chemicals
- Metallurgical industry: quality improvement of alloy steels, energy saving in metallurgical manufacturing, new alloy steel design, optimization and bottle-neck problem-solving
- Petrochemical industry: Yield and quality improvement of jet fuel, gasoline, solvent oil and others
- Automotive industry: Quality control in casting, heat treatment and electroplating processes.
- Semiconductor industry: Wafer yield control, VPTC semiconductor quality improvement.
- Materials industry: Carbon-fiber reinforced composite materials, new phosphors as lighting material
- Pharmateutic industry: Yield control in fermentation processes, drug design.
- Energy saving and environmental control: reduction of dust emission from kilns, energy saving.

REFERENCES

1. N. Y. Chen and D. Zhu, 1998, "Data mining for industrial design and diagnosis," Proceedings of the First International Conference on Multisensor-Multisource Information Fusion, Las Vegas, NV.
2. H. Liu Y. Chen and N. Y. Chen, 1994. "PCB Method applied to material design - computer aided synthesis of a super-conductor," Journal of Chemometrics, 8, 429-443.
3. P. Geladi, 1986, Analytica Chemica Acta, 185(1).
4. J. Holland, 1975, "Adaptation in nature and artificial intelligence," The U of Michigan Press, Michigan.
5. Xiaohua Hu and Nick Cercone, 1999, "Data Mining via Discretization, Generalization and Rough Set Feature Selection," Knowledge and Information Systems: An International Journal, 1(1).
6. U. Fayyad, 1996, "Knowledge discovery and Data mining: towards a unifying framework," Proceedings of KDD-96, Menlo Park, CA, AAAI Press, 82-88.
7. S. Anand, EDM, 1996, "a general framework for data mining based on evidence theory," Data and Knowledge Engineering, 18, 189-223.

Databases and Semantic Networks for Computer Design of Inorganic Materials

N.N. Kiselyova

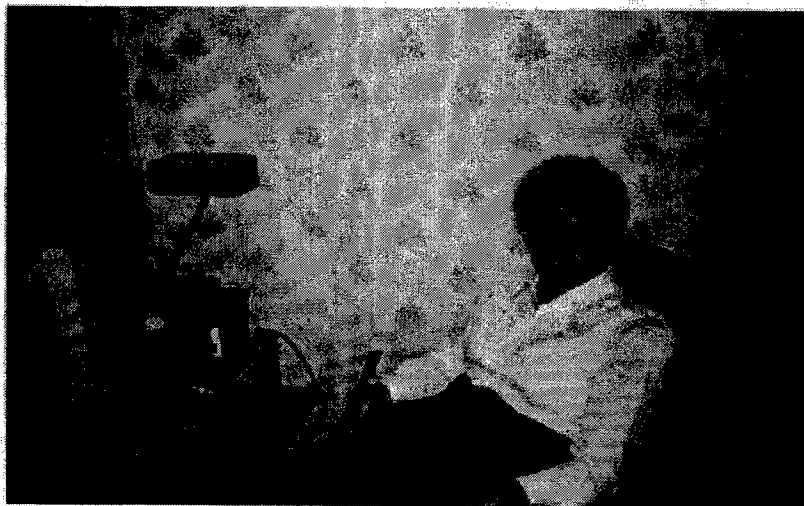
A.A.Baikov Institute of Metallurgy and Materials Science,
Russian Academy of Sciences, Moscow, Russia

Email: kis@ultra.imet.ac.ru

ABSTRACT

At present, hundreds of data bases (DBs) on substances and material properties have been developed. The prime aim of their operation is information service. Fifteen years ago, we proposed to use an extensive information of DBs not only for information service but also for searching for regularities in data and the application of these regularities for prediction of new substances. The semantic networks of special interest were used to search for regularities. Using these deduced regularities, we have predicted thousands of new compounds in ternary, quaternary and more complicated systems and have estimated some of their properties (crystal structure crystal type, melting point, homogeneity region, etc.). Comparison of our predictions with experimental data, obtained later, showed that the average reliability of predicted inorganic compounds exceeds 80 %.

DATABASES AND SEMANTIC NETWORKS FOR THE INORGANIC MATERIALS COMPUTER DESIGN

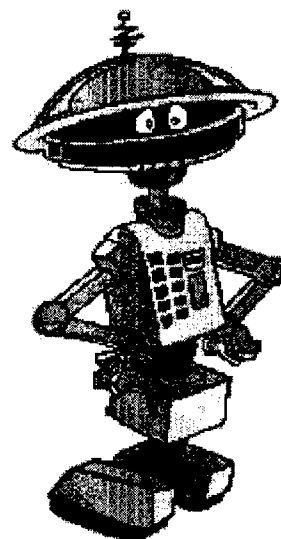


Nadezhda N. Kiselyova

**A.A. Baikov Institute of Metallurgy and Materials Science of Russian Academy of
Sciences, Moscow, Russia**

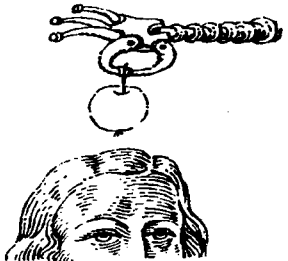
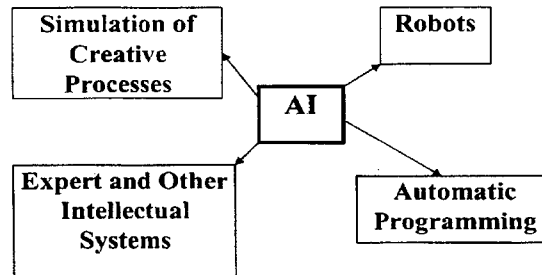
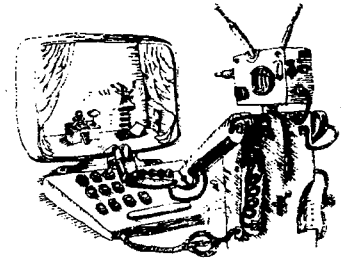
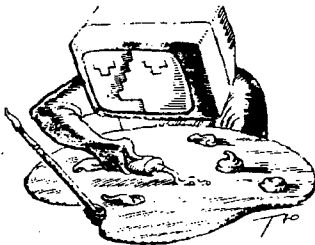
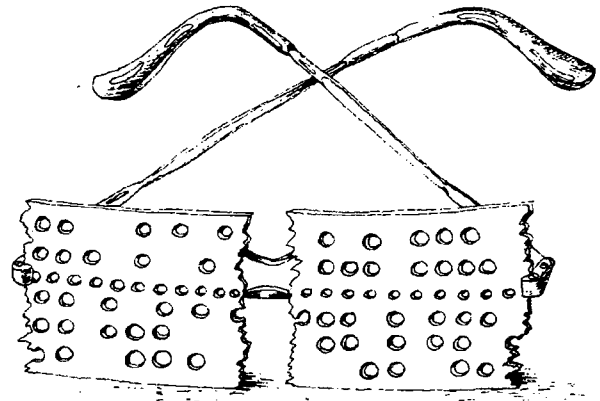
ABSTRACT

At present the hundreds of data bases (DBs) on substance and material properties are developed. The prime aim of their operation is an information service. 15 years ago we proposed to use an extensive information of DBs not only for information service but for searching for regularities in data and the application of these regularities for the prediction of new substances also. The semantic networks of special kind were used for the search for regularities. Using deduced regularities we have predicted the thousands of new compounds in ternary, quaternary and more complicated systems and estimated their some properties (crystal structure type, melting point, homogeneity region etc.). The comparison of our predictions with experimental data, obtained later, showed that the average reliability of predicted inorganic compounds exceeds 80 %.



INTRODUCTION

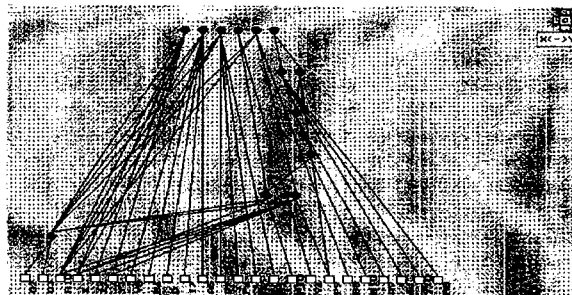
The idea of development of artificial intelligence (AI) as a model of '*homo sapiens*' occurred with the advent of "clever machines" - computers. But as years went by, ideas of AI were transformed to the development of robots performing the some of human duties under unfavourable for men conditions, to the programs extending the capabilities of dialog with computers, to the expert systems allowing the solution of tasks which defy mathematical formulation, etc. The most interesting AI applications are the data processing program systems for large symbolic information bulks. The target of this processing is a search for regularities in data. At the beginning these programs were made for robot learning in purposeful behaviour in actual practice and as a tool for an analysis of audio- and video-information to be input into computer. But it has evident soon that the area of these programs application can be much more one and they are a tool extending the human capabilities in cognition of the universe. The AI methods named as computer learning came into use widely for an analysis of geological and geophysical information with the aim of prediction of the deposits or earthquakes, for technical diagnostics for the search for machinery faultinesses, for medical diagnostics, for analysis of spectral data for the purpose of detection of various chemical and physical effects, etc. This paper presents the results of computer learning applications to data processing of a great body of information about inorganic substances with the goal of prediction of new materials with the predefined properties.



SEMANTIC NETWORKS FOR THE REPRESENTATION KNOWLEDGE ABOUT INORGANIC SUBSTANCES

The data processing is an analysis of properties and relations of objects with the aim of detection of various connections in between. Therefore the most optimal form of such a data representation in the computer memory is associative structures which allow to trace the connections without a great body of sorting information. The mathematical model of associative structures is an especial kind of graphs - the semantic networks (SNs) describing the connections between objects, their properties and status [1]. The SNs in inorganic chemistry can represent the connections between properties of components of known physical-chemical systems and properties of these systems. It is important that the chemical systems can be divide into distinguishable classes with their physical and chemical properties. These classes were determined by some concepts, for example, "systems with formation of compounds of definite composition", "compounds with definite crystal structure type", "superconductors with the nitrogen critical temperature of transition to the superconducting state", and so on. The search for sets of component properties' intervals, which cause the system membership to a certain class, is one of the chemical data processing task using the SNs. The result is a general classifying regularity (a computer form of some known concept description). This process is referred to as computer learning (or concept formation using terms of the method [1]). This method is based on a SN use. The periodicity of chemical elements' properties and the following from this fact

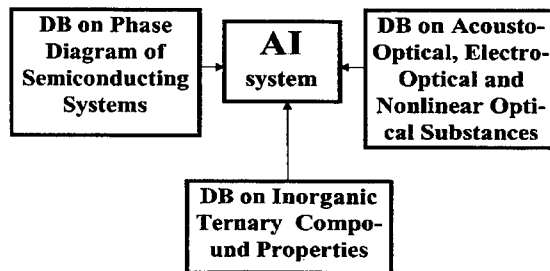
periodicity of compounds' properties, allow to use the formed computer concepts for recognition of membership of unknown chemical element set to one or other class which are described by formed concepts.



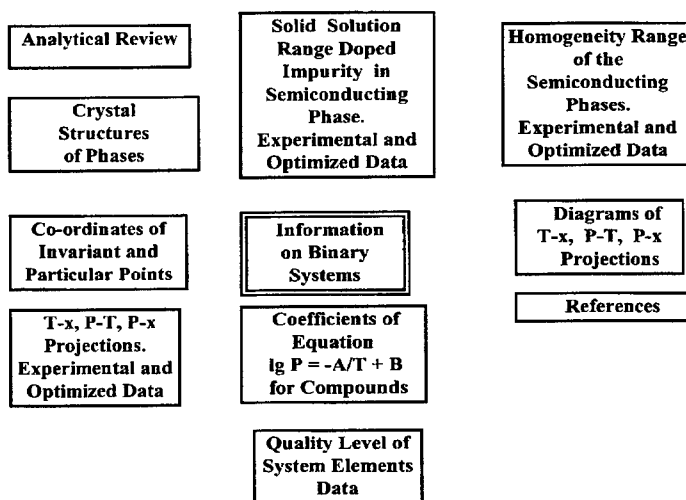
DATABASES AS A FOUNDATION OF SEMANTIC NETWORKS BUILDING

The problem of application of SNs and another empirical and semi-empirical methods for finding regularities is an use of rather complete and qualitative data. Our experience of SN applications shows that the number of erroneous predictions varies proportionally with a number of errors in experimental data to be processed and the reliability grows with an increase of initial data volume (reliability mounts to a limit with an increase of size and representativeness of learning set). Consequently, the application of the methods of the search for regularities, based on SNs, implies an use of databases, containing extensive bulks of qualitative information, as a basis. With this aim in mind we develop the DBs containing data with the qualified expert assessment. The most interesting of them are the following: DBs on materials for electronics with completely assessed information [2,3] and an inorganic ternary compound properties DB containing partially assessed information [4,5].

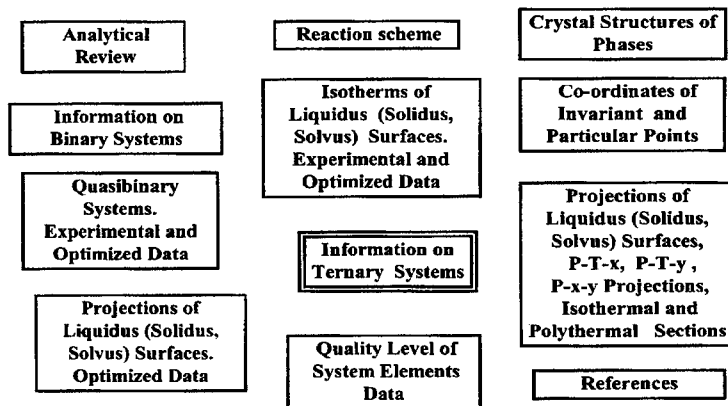
Our investigations, aimed at predicting new materials for electronics and other applications, are based on these DBs. The development of a DB is connected with building of some predefined SN that represents the objective interrelations between the properties of substances. But DB information does not provide direct answers to connection between the substance properties and constituent component properties. The AI application makes it possible to search for such connections.



1) a **phase diagram DB** of material systems with intermediate semiconducting phases [2] contains information on physical-chemical properties of the intermediate phases and the most important Pressure-Temperature-Concentration phase diagrams of semiconducting systems evaluated by qualified experts. Currently the DB contains information on several tens of semiconducting systems.

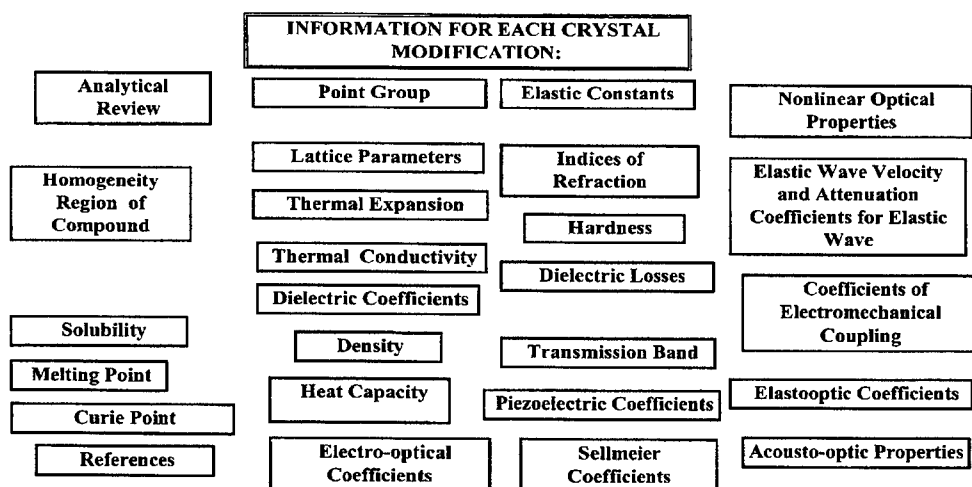


The Database of Phase Diagrams of Binary Semiconducting Systems.

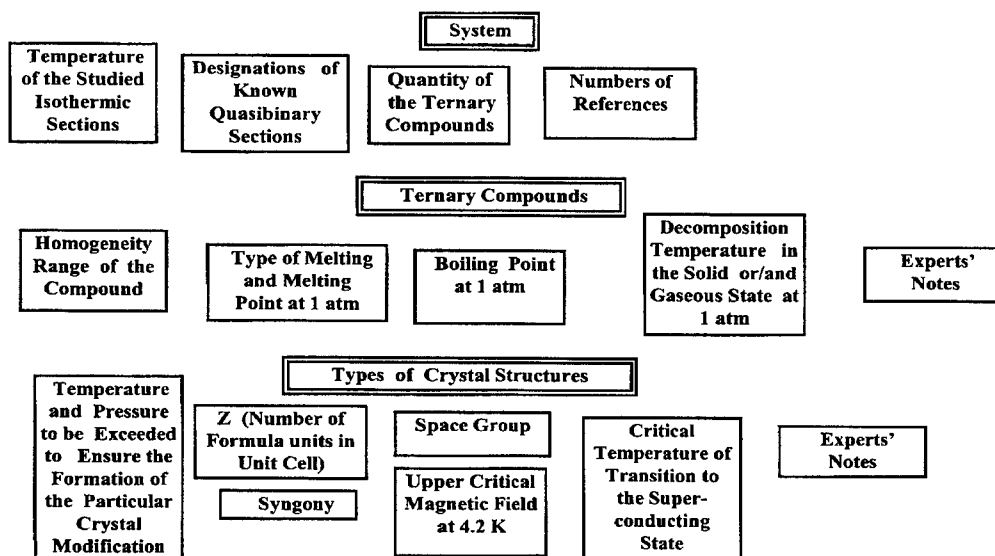


The Database of Phase Diagrams of Ternary Semiconducting Systems.

2) an acousto-, electro-, and nonlinear-optical properties [3] DB which contains information on crystals evaluated by experts. In addition, DB includes extensive graphical information about the properties of the materials



The Database on Crystals with Acousto, Electro- and Nonlinear Optical Properties "Crystal".



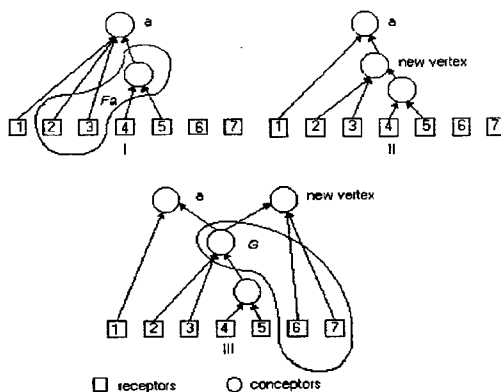
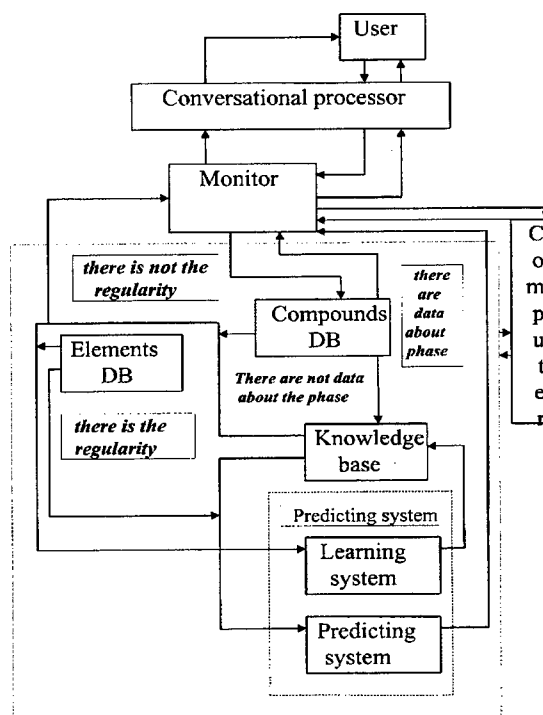
The Database of Ternary Inorganic Compounds "Phases".

3) an inorganic ternary compound properties DB was built by us in the 1970's [4,5]. It contains information about thermochemical, crystal chemical and superconducting properties on more than 37,000 ternary compounds taken from more than 11,000 publications. Some of the data has been assessed by materials experts.

APPLICATION OF ARTIFICIAL INTELLIGENCE AND DATABASES TO THE NEW INORGANIC MATERIALS COMPUTER DESIGN

It is impossible to use the DBs completely without an especial software of the search for regularities in data. During a quarter of a century we and our colleagues from Glushkov Institute of Cybernetics of National Academy of Sciences of Ukraine work at the problems of AI application to the prediction of new inorganic materials with predefined properties [6-10]. These investigations are aimed at development of the information-predicting system for inorganic materials computer design, based on SNs [6-8].

Schematic diagram of an information-predicting system



Network building [1].

vertices in built network and the choice of those ones that are the most typical for each class [1]. These vertices became the checking vertices.

We use the system that represents an initial information about known physical-chemical systems as some SNs - growing pyramidal networks (GPNs). A pyramidal network is an acyclic oriented graph having no vertices with one entering arc. If the processes of concept formation are determined in the network then the pyramidal network is designated as growing one [1]. GPN is built during the process of objects input. Each object (physical-chemical system) is put in as a set of values of the component properties with an indication of the class to which the system belongs. The nearby values of components' properties are united into one interval using an especial program or an experience of researcher. Concept formation process consists in the analysis of

(2) all properties of simple compounds - oxides, chalcogenides, halides, etc. - as required by the composition of the compounds predicted [6-10].

Our approach has made it possible to solve problems of the following types [6-10]:

- prediction of compound formation or non-formation for ternary systems;
- prediction of the possibility of forming ternary and more complicated compounds of desired composition;
- prediction of phases with defined crystal structures;
- estimation of phase quantitative properties (critical temperature of transition to superconducting state, homogeneity region, etc.).

Illustrated in Figure is the comparison between the results of predicting the compounds with composition ABO_3 [9] and the new experimental data. It is a first prediction which we carried-out 25 years ago. Only one prediction was detected to be in error ($CuGeO_3$).

A ^{II} B ^{IV}	Bc	Mg	Ca	Mn	Fe	Co	Ni	Cu	Zn	Sr	Cd	Sn	Ba	Hg	Pb
C	↔	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕
Si	+	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	+	⊕
S			⊕	⊕	⊕	⊕	+		+	⊕	⊕	⊕	⊕	⊕	⊕
Ti	↔	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	+	⊕	⊕	⊕
V			⊕	⊕	+	⊕	⊕	⊕	+	⊕	⊕	+	⊕	⊕	+
Mn	+	⊕	⊕	+	⊕	⊕	⊕	⊕	⊕	⊕	⊕	+	⊕	+	⊕
Ge	+	⊕	⊕	⊕	⊕	⊕	+	⊕	⊕	⊕	⊕	⊕	⊕	+	⊕
Se		⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	⊕	+	⊕	⊕	⊕
Zr	↔	↔	⊕	+	↔	+	+	+	⊕	⊕	⊕	+	⊕	+	⊕
Mo		⊕	⊕	⊕	⊕	⊕	⊕	+	⊕	⊕	+	+	⊕	+	+
Tc			⊕	+	+	+	+	+	+	⊕	+	+	⊕	+	⊕
Ru			⊕	+	+	+	+	+	+	⊕	+	+	⊕	+	+
Sn			⊕	⊕	+	⊕	+	+	⊕	⊕	+	+	⊕	+	⊕
Te		⊕	⊕	⊕	+	⊕	⊕	⊕	⊕	⊕	⊕	+	⊕	⊕	⊕
Ce	↔	⊕	⊕					⊕	+	⊕	⊕		⊕	+	⊕
Pr								+	+	⊕	+		⊕	+	+
Tb	-	-	+					+	+	⊕	+		⊕		
Hf		↔	⊕	+	+	+	+	+	+	⊕	⊕	⊕	⊕		⊕
Ta			+	+	+	+	+	⊕	+	+	+	⊕	⊕	+	
Os	-	-	⊕	+	+	+	+	+	+	⊕	⊕	+	⊕		
Ir			⊕	+	+	+	+	+	+	⊕	+	+	⊕		
Pb			⊕	⊕	⊕	+	+	⊕	⊕	⊕	⊕	⊕	⊕	⊕	
Po			+	+	+	+	+	+	+	⊕	+	+	⊕		
Th	↔	⊕	⊕	+	+	+	+	+	+	⊕	⊕	+	⊕		⊕
U	↔	↔	⊕	+	+	⊕	+	+	+	⊕	⊕	+	⊕		

Designations: + - formation of compound with composition ABO_3 is predicted; - - formation of compound with composition ABO_3 is not predicted; ⊕ - compound with composition ABO_3 was synthesized and appropriate information was used in the computer learning process; ↔ - compound with composition ABO_3 does not exist under normal conditions and this information was used in the computer learning process; ⊙ - predicted formation of compound with composition ABO_3 which was confirmed by experiment; ⊗ - predicted formation of compound with composition ABO_3 which was not confirmed by experiment; empty square - indeterminate result.

Part of a Table Illustrating the Prediction of Compounds with the Composition ABO_3 .

Compounds/Systems	Characteristics to be Predicted	Experimental Tests for March 1999	Error of Prediction, %
ABX (X=Se,Te)	Compound formation	99	39
ABX ₂ (X=O,S,Se,Te)	Compound formation	328	10
ABX ₃ (X=O,F,S, Cl,Se,Br,Te,I)	Compound formation	381	14
ABX ₄ (X=O,F,Cl,Br,I)	Compound formation	393	5
A ₂ BX ₂ (X=S,Se)	Compound formation	24	9
AB ₂ X ₄ (X=O,F,S, Cl,Se,Br,Te,I)	Compound formation	746	15
A ₂ B ₂ X ₂ (X=O,S,Se)	Compound formation	97	26
A(Hal) ₂ - B(Hal)	Systems w/ compounds	108	10
AB ₂ X ₄ (X=O,S,Se,Te)	Structure type	367	6
ABX (X=Al,Si,P,Ga,Ge,As,Pd,In,Bi)	Structure type	46	50
ABO ₃	Perovskite structure	186	13
A ₂ B ₂ O ₇	Pyrochlore structure	73	18
AB ₂ X ₂ (X=Al,Si,P, Ge,As, Sb)	ThCr ₂ Si ₂ structure	157	6
ABX ₂ (X=Al,Co,Ni, Cu,Ga,Pd,In)	MnCu ₂ Al structure	55	13
A ₂ (SO ₄) ₃ * B ₂ (SO ₄) ₃ and A(NO ₃) ₃ * B(NO ₃) ₃	Compound formation 1:1	130	4
ABDO ₄	Compound formation	22	6
		Average	= 15 %

Using AI approach described above we have predicted the formation of thousands of new compounds in ternary, quaternary and more complicated systems. These compounds were then searched for new magnets, semiconductors, superconductors, electro-optical, acousto-optical, nonlinear optical and other materials required for new technologies [6-10]. The comparison of these predictions with the experimental data, obtained later, showed that the average reliability of predicted compounds exceeds 80 %.

REFERENCES

1. V.P.Gladun, *Processes of Formation of New Knowledge*, 1994. SD "Pedagog", Sofia, 192 p.
2. V.S.Zemskov, N.N.Kiselyova, N.N.Kiselyova, et al., 1995. DIAGRAMMA Database on Phase Diagrams of Semiconductor Systems. *Inorganic Materials*, 31(9) 1096-1100.
3. N.V.Yudina, V.V.Petukhov, E.A.Chermushkin, et al., 1996. Data Bank on Acoustooptical, Electrooptical, and Nonlinear Optical Properties of Materials. *Crystallography Reports*, 41(3), 464-468.
4. E.M.Savitskii, N.N.Kiselyova, B.N.Pishik, et al., 1984. Development of Data Bank on Ternary Inorganic Phase Properties. *Doklady AN SSSR*, 279(3), 627-629.
5. N.N.Kiselyova, N.V.Kravchenko, V.V.Petukhov, 1996. Database System on the Properties of Ternary Inorganic Compounds (IBM PC Version). *Inorganic Materials*, 32(5), 567-570.
6. N.N.Kiselyova, 1997. Application of Artificial Intelligence Methods to Inorganic Compounds Computer Design. *Perspektivnye Materialy*, 4, 5-21.
7. E.M.Savitskii, V.B.Gribulya, N.N.Kiselyova, et al., 1990. *Prediction in Material Science Using Computer*, Nauka, Moscow, 86 p.
8. N.N.Kiselyova, 1993. Information-predicting systems for the design of new materials. *J.Alloys and Compounds*, 197(2), 159-165.
9. N.N.Kiselyova, B.I.Pokrovskii, L.N.Komissarova, 1977. Simulation of forming complicated oxides from initial components using the cybernetic method of concept formation. *Zh. Neorgan. Khimii*, 22(4), 883-886.
10. N.N.Kiselyova, 1998. Computational materials design using artificial intelligence methods. *J.Alloys and Compounds*, 279(1), 8-13.

First-Principles Calculations for Materials Science: Their Power and Limitations

Wanda Andreoni

IBM Research Division, Zurich Research Laboratory,
CH-8803 Rueschlikon, Zurich, Switzerland

The use of parallel computers and that of sophisticated computational algorithms has made first-principle calculations feasible also for realistic models of systems of technological interest. Several successes have been obtained recently. In spite of this, there are still a number of problems to be solved, before this type of approach can assume a leading role in the investigation of materials. This talk will report on recent applications of parameter-free molecular dynamics to a variety of systems, with emphasis on the method, on the comparison of the results with experiment, on its useful outputs and also on its current limitations.

Interplay between Large Materials Databases, Semi-Empirical Approaches, Neuro-Computing and First Principle Calculations

Pierre Villars*, Steven R. LeClair** and Shuichi Iwata***

* Material Phases Data System (MPDS), CH-6354 Vitznau, Switzerland

** U.S. Air Force Research Laboratory,

2977 P Street, Suite 13, Wright-Patterson AFB OH 45433-7746, USA

*** RACE, Faculty of Engineering, The University of Tokyo,

7-3-1 Hongo, Tokyo 113, Japan

Materials design is still mainly based on known concepts in materials science and intuition of the experimentalists. Analyzing the conditions that make it possible to search for materials science concepts, shows that it was not a new technique, a unique experimental observation, or an abstruse theory, which formed the take-off point. It was rather the amassing of a critical volume of experimentally-determined data in the literature that permitted an individual with deep insight to perceive an underlying pattern not previously apparent. Extending these facts to a new area of materials design leads to the following four key-points:

1. Creation and use of huge, critically evaluated materials databases which comprehensively covers the published world literature (materials databases).
2. Computer-aided reduction of elemental parameters and systematic combinations of them to find the relevant feature sets which can link materials properties qualitatively with the chemical species present (semi-empirical approaches).
3. Refinement and optimisation of qualitatively-obtained results under (2) with the help of neuro-computing leading to more explicit quantitative results.
4. Focusing on predicted, most-promising materials systems with the aim to reduce the experimental work for verification, as well as trying to create a theoretically-based explanation for such quantitative results (first-principle calculations).

Materials Databases:

The amount of critically-evaluated materials data has reached a respectful level, but it is still far from comprehensive. Below are the 6 most significant materials databases available in electronic form:

ICSD - This Inorganic Structure Database is maintained by the Fachinformationszentrum in Karlsruhe, Germany and contains crystallographic data for inorganic compounds.

CRYSTMET - This Intermetallic Structure Database was maintained until April 1, 1997 by the National Research Council of Canada NRCC. CISTI, hard copy versions are Pearson's Handbook of Crystallographic Data for Intermetallic Phases, ASM International 1991 and Pearson's Desk Edition, ASM International, 1997 which contain crystallographic data for intermetallics and alloys.

ICDD PDF2 - This Powder Diffraction Patterns Database is maintained by the International Centre for Diffraction Data in Swarthmore, Philadelphia which contains mainly measured powder patterns.

**BINARY ALLOY -
PHASE DIAGRAMS** - This CD-ROM is maintained by ASM International (Editor-in-Chief: T.B. Massalski)

INTERPLAY BETWEEN LARGE MATERIALS DATABASES, SEMI-EMPIRICAL APPROACHES, NEURO-COMPUTING and FIRST PRINCIPLE CALCULATIONS

by

P. Villars, Materials Phases Data System (MPDS), CH-6354 Vitznau, Switzerland

S. LeClair, Wright-Patterson AFB, 2977 P Street, Suite 13, Dayton OH 45433-7746, USA

S. Iwata, RACE, Faculty of Engineering, University of Tokyo, 7-3-1 Hongo, Tokyo 113, Japan

M. Berndt, Crystal Impact, Postfach 1251, D-53002 Bonn, Germany

K. Brandenburg, Crystal Impact, Postfach 1251, D-53002 Bonn, Germany

**General Idea of 'INTERPLAY' with the ultimate aim of
'Virtual Materials Design' capability**



**New Materials - from trial and error - to combinatorial
chemistry -to pattern directed discovery - to first principle
calculations**

Materials design is still mainly based on the in materials science known concepts and intuition of the experimentalists. Analyzing the conditions that make it possible to search for the in materials science known concepts shows that it was not a new theory, a unique experimental observation, or an abstruse theory which formed the take-off point. It was rather the amassing of a critical volume of experimentally determined data in the literature that permitted an individual with deep insight to perceive an underlying pattern not previously apparent.

Extending these facts to a new area of materials design leads to the following four key-points:

- ① The creation and the use of huge, critically evaluated materials databases which comprehensively covers the published world literature.

→→→ (large electronic materials databases)

- ② Computer-aided reduction of the elemental property parameters and systematic combinations of them to find the relevant 3D-feature sets which qualitatively can link materials properties with the chemical species present

→→→ (semi-empirical approaches)

- ③ Refinement and optimization of the qualitatively obtained results under
 - ② with the help of neuro-computing leading to quantitative results.

→→→ (neuro-computing)

- ④ Focusing on predicted, most promising materials systems with the aim to reduce the experimental work for its verification, as well as trying to create a theoretical based explanation for such quantitative correlations between materials properties and its constituent elemental property parameters.

→→→ (first principle calculations)

Large electronic Materials Databases

The amount of critically evaluated materials data starts to reach an acceptable magnitude,

but is still far away from calling it comprehensive,

below are the 7 most significant electronic materials databases listed:

★ ICSD

This Inorganic Structure Database is maintained by the Fachinformationszentrum in Karlsruhe, Germany and contains crystallographic data for Inorganic Compounds.

★ Pearson's Handbook *(by P. Villars and L.D. Calvert)*

Electronic version available by MPDS, its hardcopy versions are Pearson's Handbook of Crystallographic Data for Intermetallic Phases, ASM International, 1991 and Pearson's Desk Edition, ASM International, 1997 which contain crystallographic data for Intermetallics and Alloys.

★ ICDD PDF2

This Powder Diffraction Patterns Database is maintained by the International Centre for Diffraction Data in Swarthmore, Philadelphia which contains mainly measured powder patterns and most recently calculated ones (from ICSD data).

★ Binary Alloy Phase Diagrams CD-ROM *(by T. Massalski, H. Okamoto)*

This CD-ROM is maintained by ASM International.

★ Ternary Alloy Phase Diagrams CD-ROM *(by P. Villars, A. Prince, H. Okamoto)*

This CD-ROM is maintained by ASM International.

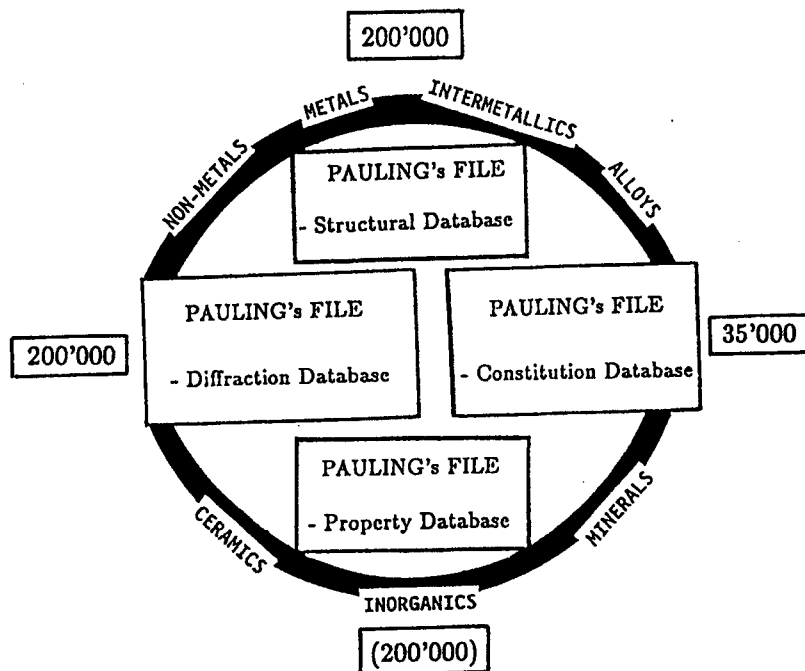
★ Landolt-Boernstein CD-ROM

Recently a series of the many-volume handbooks are available as CD-ROM, Springer-Verlag.

★ LPF, Linus Pauling File *(by editor-in-chief P. Villars)*

Is a Basic Database for Alloys, Intermetallics and Inorganics. This file is now in the process to be build up by the Japan Science and Technology Corporation JST in Japan and MPDS in Switzerland and is planned to enter the yearly update stage in 2007 which covers crystallographic data, powder patterns, intrinsic physical property data and phase diagrams (see plot 1).

CONCEPTS of the PAULING's FILE Project



+

PAULING's
Retrieval/Visualization Software Tools

||

PAULING's Materials Design System

- fully relational database
- data fully standardized
- data fully consistent
- data critically evaluated
- very comprehensive world literature coverage
 - long term approach

Semi-empirical approaches

There exists in the world literature a whole range of 'highest quality' correlations between materials properties and the chemical species present. To all of them is common that they were found by semi-empirical approaches based on a small to large amount of experimentally known data.

A comprehensive review is given in the book 'Intermetallic Compounds, Principles and Practice (Volume 1), edited by J.H. Westbrook and R.L. Fleischer, John Wiley & Sons (1995).

Here we show the first results of an automated 'discovery tool' to search systematically for the relevant 3D-feature sets (derived from elemental property parameters of the chemical species present) to correlate qualitatively materials properties with the chemical species present.

This is done in three major steps:

- ① Collection of all published elemental property parameters and find the most independent parameter sets within them

We found in total over 300 parameter sets, of which for 42 sets complete parameter sets are published

These can be grouped into the following:

*7(8) groups of elemental property parameter sets, here also called the
7(8) "Factors":*

(for each elemental property parameter set we choose the most accurate one)

Atomic weight Factor (SI-unit: kg) weight atomic or atomic number

Electro-chemical Factor (SI-unit: mol) electronegativity after M&B

Frequency Factor (SI-unit: s) magnetic resonance

Heat Factor (SI-unit: K) temperature melting

Size Factor (SI-unit: m) radii pseudo-potential after Zunger

Valence electron Factor (SI-unit: A) valence electron number

[Chemistry Factor: number of electrons + s,p,d,f shell electrons after Thaler]

Optical Factor (SI-unit: cd) ?

② Building an automatic generator for 3D-feature sets resulting from combinations of **elemental property parameters** and **mathematical operations**.

The number of 3D-feature sets has to be chosen very carefully because otherwise the total number of 3D-feature sets becomes 'astronomic'. To start e.g. with 6 elemental property parameters and 5 basic mathematical operations (sum, difference, ratio, product, maximum) results in 30 combinations. Three of them gives one 3D-feature set to be investigated, that means $(30 \cdot 29 \cdot 28) / (3 \cdot 2) = 4060$ 3D-feature sets.

e.g. 1 additional elemental property parameter would give 6545 3D-featuresets, and 1 additional operation would result in 7140 3D-feature sets, and 1 additional elemental property parameter + 1 additional operation results in 11480 3D-feature sets.

③ Automatic High-Quality Separation Detection (+ its Visualization).

What humans normally do using their eyes and brains is: finding some kind of border between areas and counting how many points with the same materials property are the same at the same side of the border. In three dimensions one would have to find a more or less complex surface, and though algorithms to find them exists, they would take too much calculation time to be applied on thousands of 3D-feature sets. So the detection algorithm must work on a much simpler basis.

Our starting point is: "If a separation is 'good' then many points' nearest neighbour(s) must have the same materials property. And this can be detected with simple distance calculations. In principle all distances from each point have to be calculated to find out, what is the shortest and therefore what is the nearest neighbour.

So what about the 'other' neighbours ? A really good separation means that points of the same materials property give clusters, which are as big as possible. Therefore we investigate the best results received by looking just at the nearest neighbour again in more detailed way, by increasing the numbers of neighbours to e.g. 50 and follow its separation behaviour (see 'hits' vs. number of neighbours, plot 2).

This can be best demonstrated on a practical example.

☉ **Considering the following 6 elemental property parameters:**

- 1) atomic number
- 2) electronegativity after M&B
- 3) magnetic resonance
- 4) temperaure melting
- 5) radii pseudo-potential after Zunger
- 6) valence electron number

☉ **Considering the following 5 mathematical operations:**

- 11) sum
- 12) difference
- 13) ratio
- 14) product
- 15) maximum

☉ **Considering 6431 ternary systems**

4104 formers +
(from Pearson's Handbook)

2327 nonformers

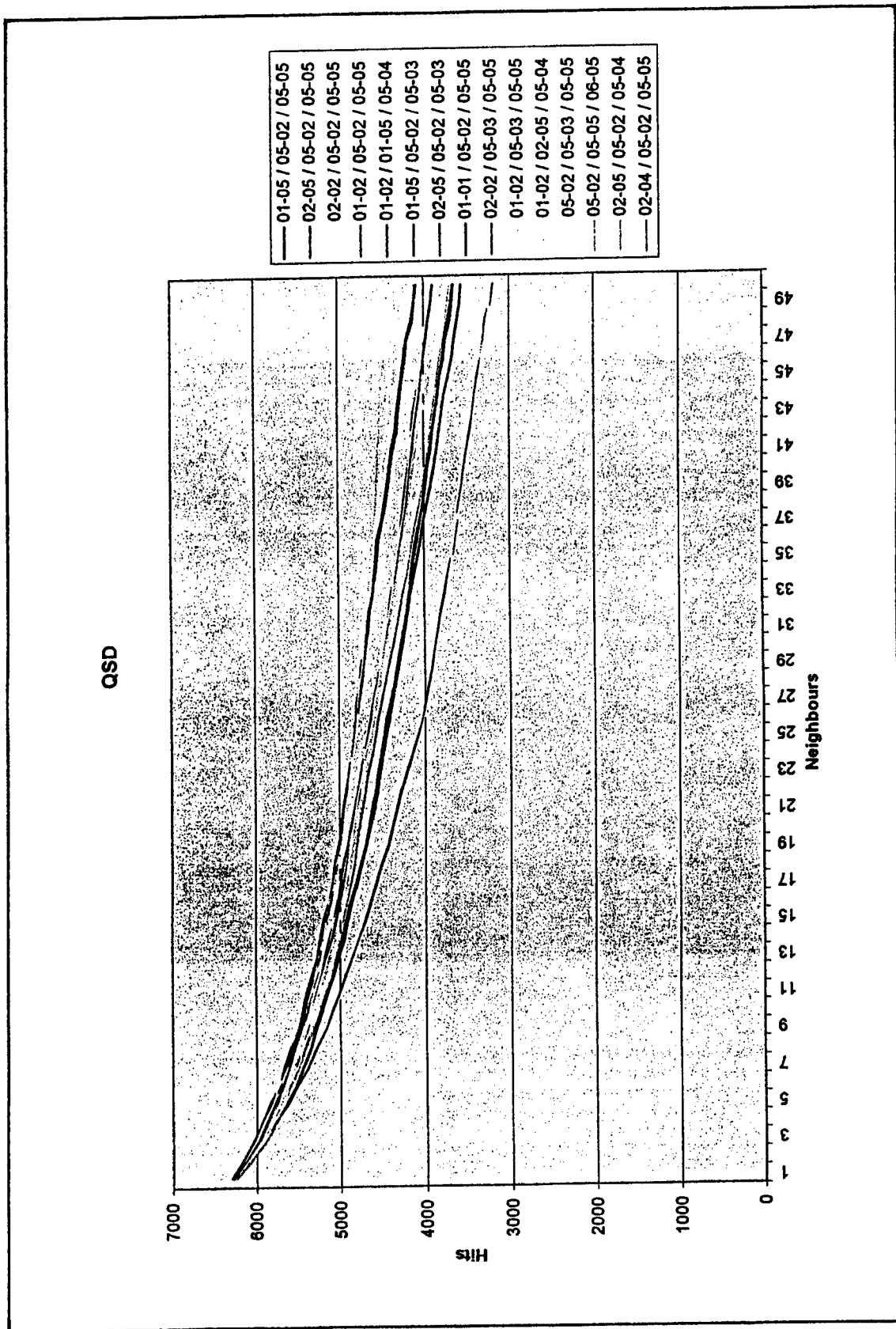
(derived from 676 binary nonformers, from Binary Alloy Phase Diagrams CD-ROM)

Below are listed the results for the 20 'best' from all

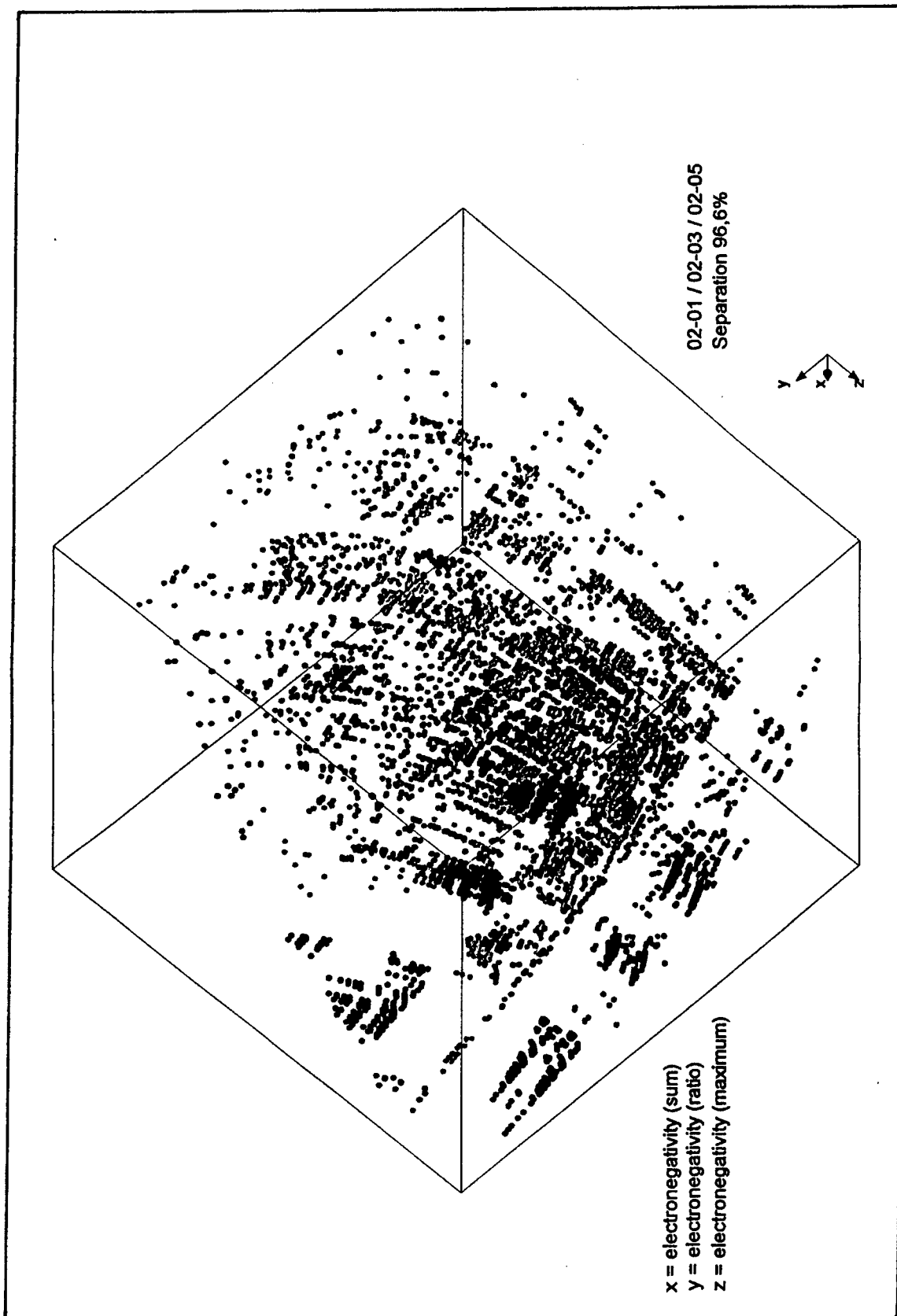
$(30 \times 29 \times 28) / (3 \times 2) = 4060$ 3D-feature sets (see plots 3 + 4):

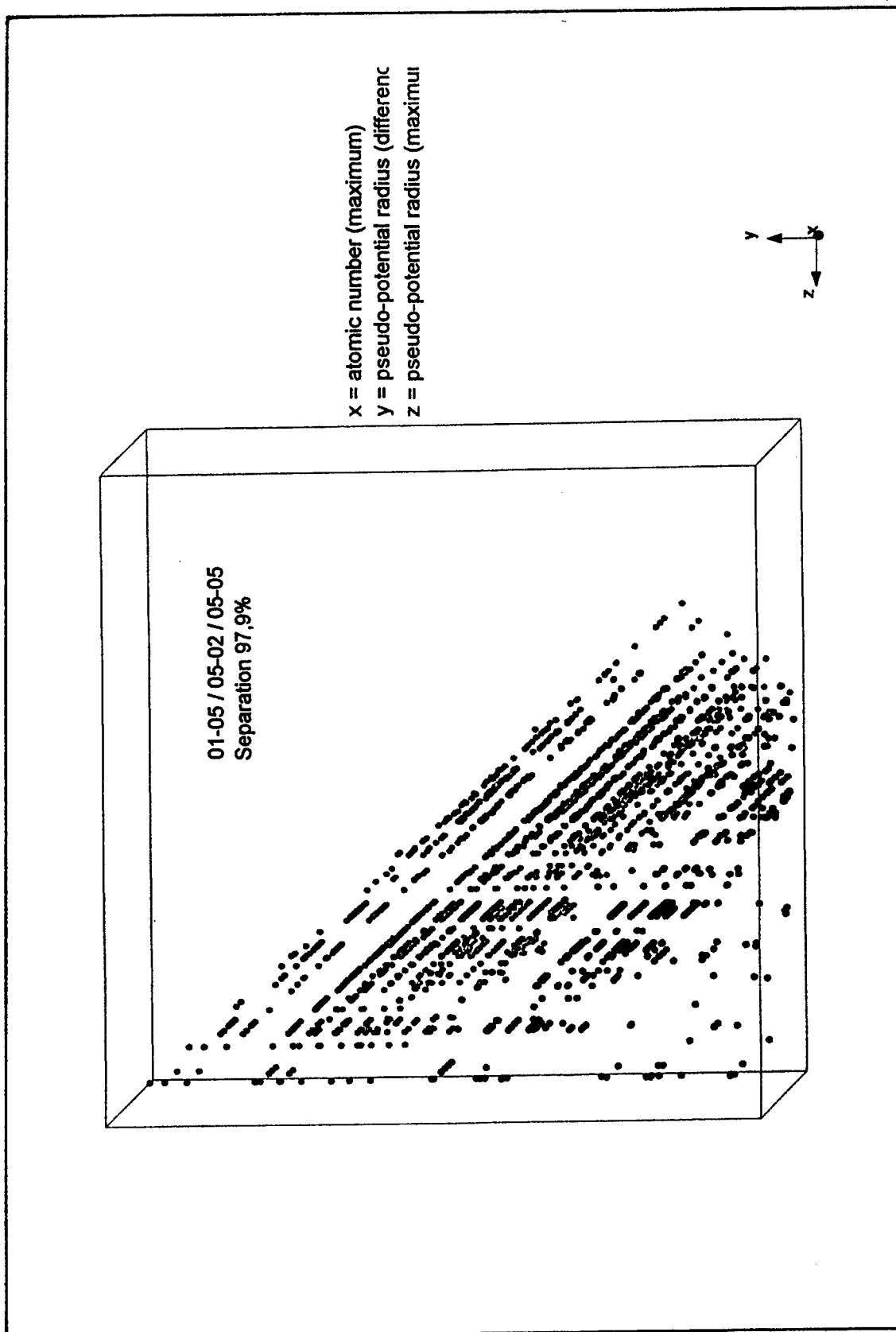
This means e.g. for 6431 ternary formers/nonformers infos $6431 \times 6430 / 2 = 20'675'665$ distances for each 3D-feature set.

1	01-15 / 05-12 / 05-15	6297 (97.92%)
2	02-15 / 05-12 / 05-15	6288 (97.78%)
3	02-12 / 05-12 / 05-15	6282 (97.68%)
4	01-12 / 05-12 / 05-15	6276 (97.59%)
5	01-12 / 01-15 / 05-14	6268 (97.47%)
6	01-15 / 05-12 / 05-13	6261 (97.36%)
7	02-15 / 05-12 / 05-13	6257 (97.29%)
8	01-11 / 05-12 / 05-15	6256 (97.28%)
9	02-12 / 05-13 / 05-15	6250 (97.19%)
10	01-12 / 05-13 / 05-15	6249 (97.17%)
11	01-12 / 02-15 / 05-14	6248 (97.15%)
12	05-12 / 05-15 / 06-15	6245 (97.11%)
13	05-12 / 05-13 / 05-15	6245 (97.11%)
14	02-15 / 05-12 / 05-14	6242 (97.06%)
15	02-11 / 05-12 / 05-15	6241 (97.05%)
16	02-14 / 05-12 / 05-15	6241 (97.05%)
17	01-15 / 05-12 / 05-14	6239 (97.01%)
18	01-15 / 05-13 / 05-14	6236 (96.97%)
19	02-12 / 05-12 / 05-13	6236 (96.97%)
20	01-15 / 05-13 / 05-15	6235 (96.95%)



plot 2





plot 4

Neuro-computing

After such qualitative correlations have been discovered the following neuro-computing approaches showed to be

striking in improving the qualitative correlations clearly towards quantitative correlations

The following approaches showed to be very successful:

- Function Approximation - Ensemble Approach
- Function Approximation - Orthogonal Approach
- Function Approximation - Auto-Associative Filtering
- Clustering and Visualization

In this session you will hear several contributions showing the power of neuro-computing.

Therefore in short the processing sequence of our Materials Design capability is: ①②③④

- ① preparing a 'highest quality data-set' from materials property of interest using large electronic materials databases. It is very important to invest some time to evaluate the data as much as possible
(otherwise garbage in → garbage out)

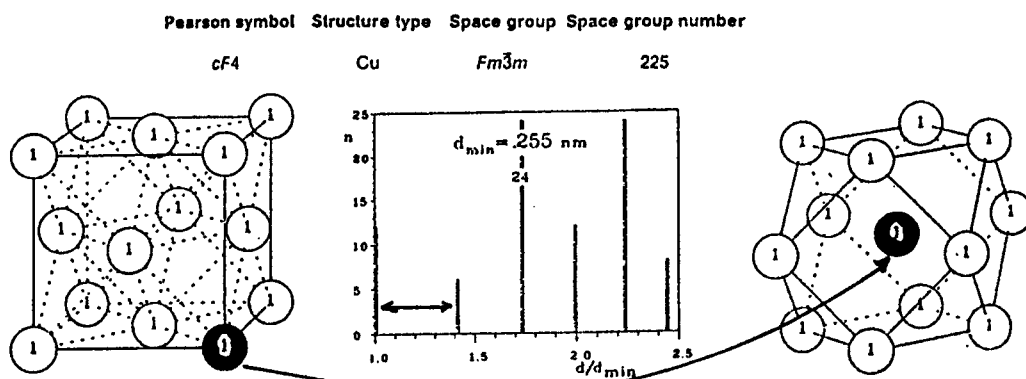
- ② using elemental property parameters + mathematical operations for the generation of 3D-feature sets

(investigating the influence of exchanging different elemental property parameters belonging to the same factor, e.g. size factor: pseudo-potential, ionic, covalent, metallic radii)

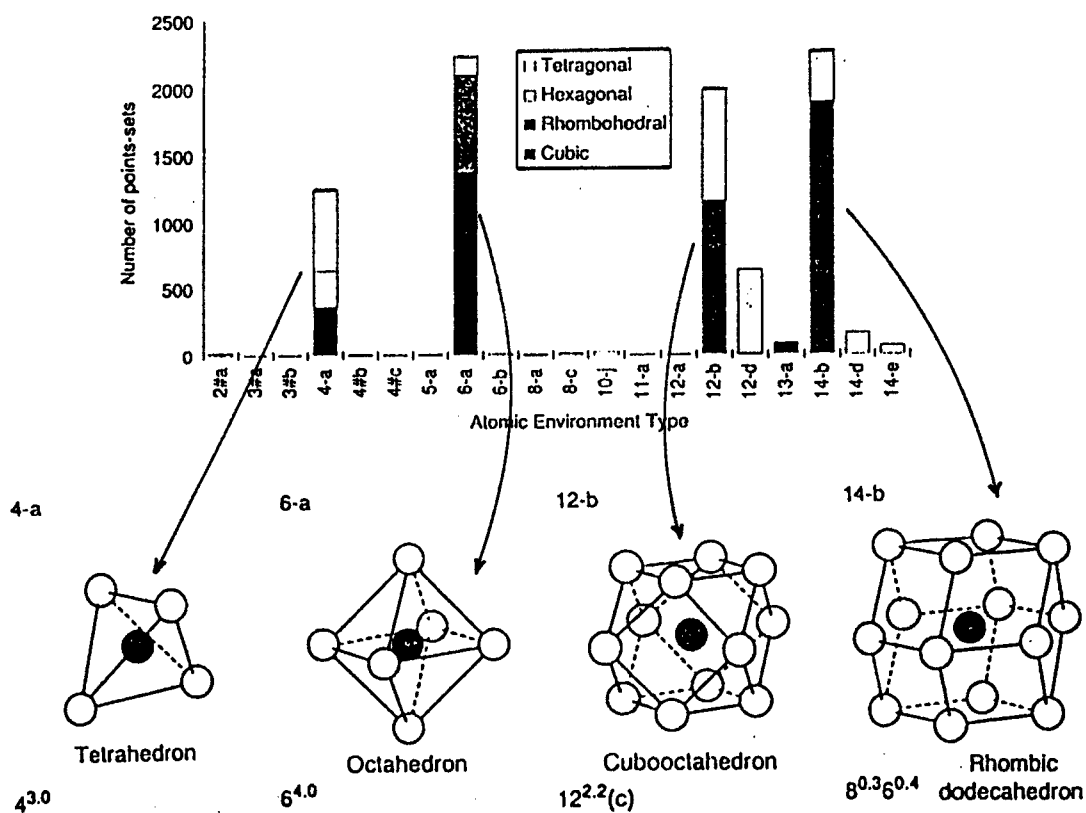
- ③ testing which 3D-feature sets are 'best' for the individual materials design problems, the 'best' solutions have to be investigated in close interaction with the results in neuro-computing. With that we have a tool to complement the 'intuition of the scientist' by a very robust method

- ④ applying different neuro-computing approaches and by optimization of the different results predicting materials properties with highest accuracy rates
(e.g. in case of ternary formers/nonformers over 99% based on experimental infos for about 6000

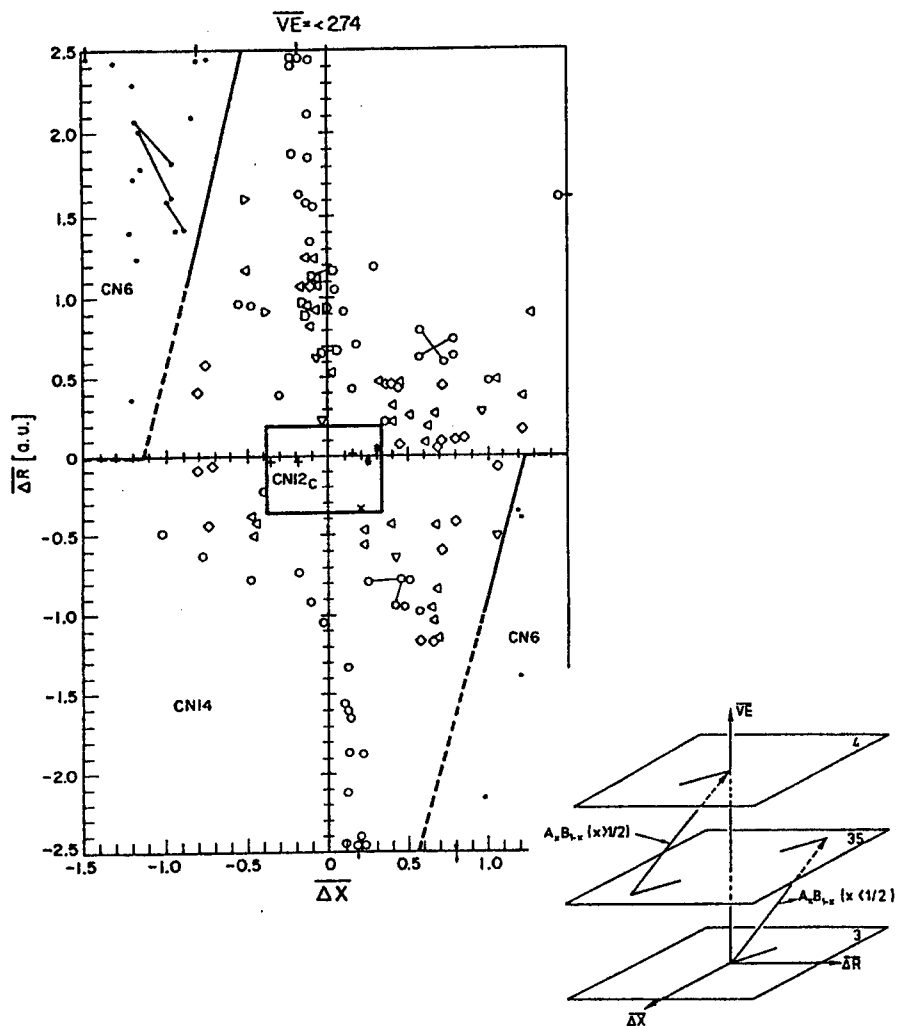
Single Environment Structure Types



Most frequently occurring AETs in the Single Coordination Types



Semi-empirical approach (*measured data*)



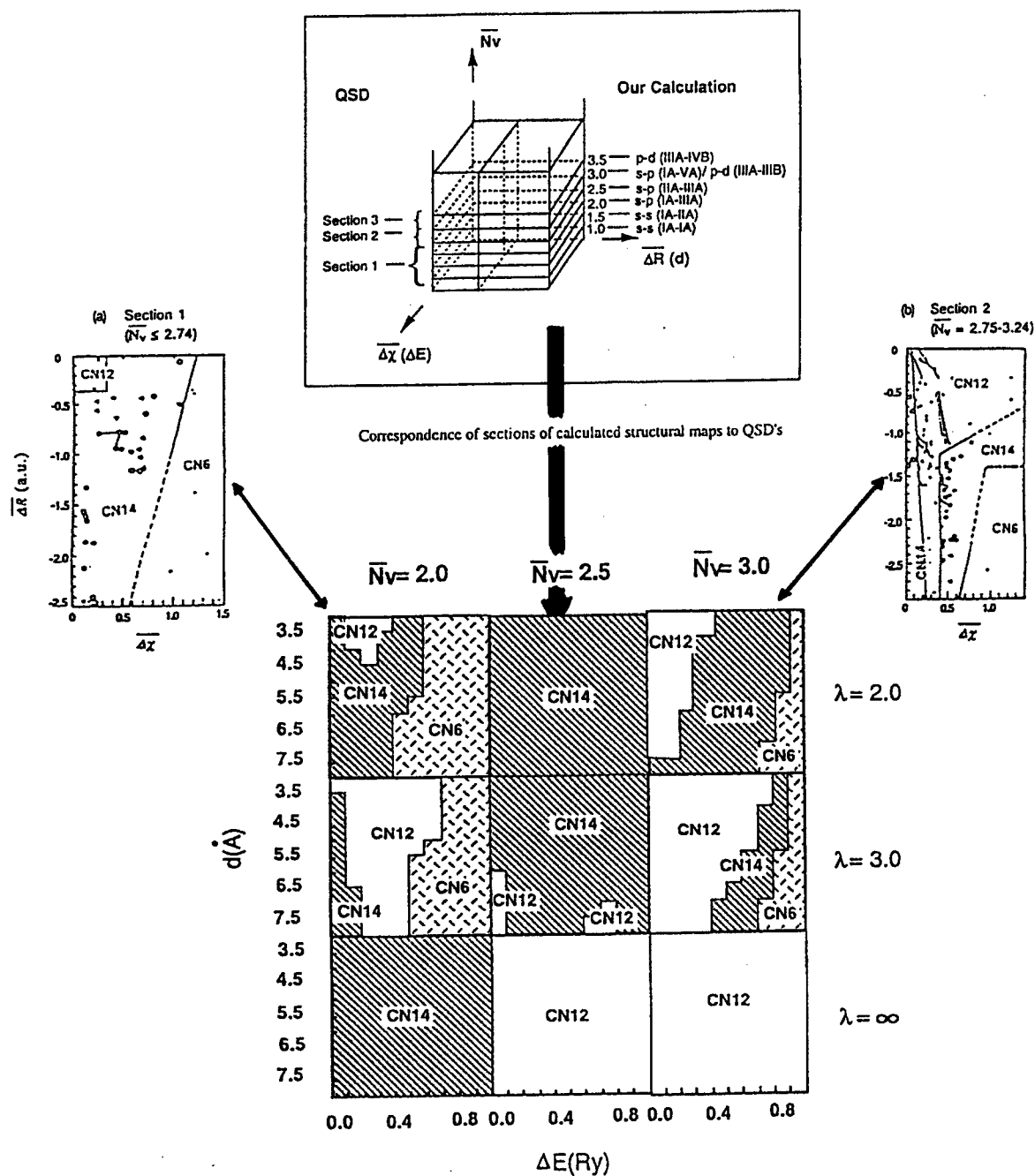
Atomic Property Expressions (Variables)

$$x = \overline{\delta R}_{s+p}^z = \delta R_{s+p}^z(\text{conc})$$

$$y = \overline{\delta X}^{MB} = \delta X^{MB}(\text{conc})$$

$$z = \overline{V} = V(\text{conc})$$

First principle calculation (*calculated data*)



Calculated structural map of s-p bonded AB compounds using the difference of valence electron orbital energies of atoms and bond length

First principle calculations

Even after establishing above mentioned quantitative correlations the number of to be experimentally verified predictions is in general too high because most advanced materials of today's interest are ternary or quaternary systems. In the case of the ternary systems there exists 161'700 and for the quaternary systems 3'921'225 potential systems. In addition, in order to establish structures and phase relationships, one has to prepare and investigate at least 10 times more samples per system by going e.g. from ternary to quaternary.

In the best case from first principle calculations one can expect to create a theoretical based explanation for such quantitative correlations between materials properties and elemental physical parameters. With that we will get a 'handle' on the processing control of its 'materials production (with pre-defined property)'

Several years ago the Structure Stability Maps, also called Quantum Structure Diagrams (QSD) was found by us using a semi-empirical approach looking at the simplest structures the single-environment compounds. Each atom within such a compound has geometrical the same atomic environment (coordination polyhedra), see plots 5 + 6.

The major achievement of the QSD was to correlate 'hidden' structural facts with elemental property parameters (3D-feature set) and achieving separation between materials with different AETs with accuracies of > 98%.

Recently Ying Chen et al. presented a simple model for studying the structure stability of atomic environments of AB intermetallic compounds. The relative stability of the four most common atomic environment types (AETs) have been systematically calculated within a tight-binding model.

The calculated three-dimensional QSD uses:

- the difference of the valence electron orbital energy of an atom, ΔE
 - the distance between atoms, d
 - the average number of electron per atoms N_v
- (+ with constant power indexes ' λ ' of the repulsive potential term)

compared with the semi-empirical found QSD uses:

- the concentration-weighted electronegativity difference, $\Delta E(M\&B)$
- the concentration weighted pseudo-potential radii difference, $\Delta R(Z)$

In summary to achieve our ultimate goal of having an accurate 'Virtual Materials Design' capability

we need:

INTERACTIONS

- ① access to huge materials databases
(high-quality data)
- ② method to systematically discover the relevant correlation between elemental property parameters (3D-feature sets) and materials property
(semi-empirical approach)
- ③ Neuro-computing to optimize results
and make highest accuracy predictions >99%
(neuro-computing)
- ④ First principal calculations to create a theoretical based explanation for such quantitative results, to get a 'handle' on the processing control of its 'production of the wanted materials'
(first principle calculation)

PREDICTION STEPS

*Being able to predict materials properties
one has to obey the following sequence:*

- ① formers versus nonformers
- ② stable compositions within the formers
- ③ crystal structure of the stable compositions within the formers
- ④ correlate materials properties with its crystal structure

WATCH

① Do not miss any elemental property parameters (3D-feature sets)

② Minimize the number of prediction steps
by covering large materials groups

(e.g. treat binay, ternary, quaternary systems together instead just binaries at a fixed stoichiometry, this means increase data amount as much as possible!)

③ Maximize the accuracy of the prediction of each prediction step to
over > 99 %

(overall accuracy is most important)

④ Reduce the number of crystal structures by grouping structures having
the same gross-feature, by moving from the 'classical' space group
description to the atomic environment description

(at present about 6'000 structure types have been published)

⑤ Correlate materials properties with its crystal structures
(be aware that this is a necessary but not sufficient condition)
using the atomic environment description

(e.g. most materials properties are found in less than 24 'classical' different crystal structures)

Software Package "Materials Designer" and its Application in Materials Research

Chen Nianyi^{*}, Lu Wencong^{}, Chen Ruiliang^{*}, Qin Pei^{*}**

^{*}Shanghai Institute of Metallurgy, Chinese Academy of Sciences, Shanghai, China

^{**}Department of Chemistry, Shanghai University, Shanghai, China

ABSTRACT

A useful software, "Materials Designer", has been built based on a series of new computation methods. It can be used to help scientists or engineers to solve a series of problems in research work, including the following purposes: (1) to optimize the technological conditions in materials preparation, in order to achieve better performance of the products; (2) to predict the properties and phase components of unknown materials. Examples of applications are described to show the usefulness of this software.

INTRODUCTION

There are two problems of general significance in materials research:

- (1) How to find the structure-property relationship of materials, in order to predict the physico-chemical properties of unknown materials systems;
- (2) How to find the best conditions of preparation of materials, in order to make experimental design for new materials preparation. The software package "Materials Designer" has been used in our laboratory for solving some topics related to the above-mentioned problems, with good results.

COMPUTATIONAL METHODS

The software "Materials Designer" is a comprehensive system consisting of several modules, including a series of pattern recognition methods: PCA, PLS, Fisher method, LMAP, MREC and ENVELOP methods. ANN and nonlinear regression modules are also included in this software. Here LMAP, MREC and ENVELOP methods are developed by ourselves. Combining with the HIDDEN PROJECTION method developed by ourselves, software "Materials Designer" provides a very powerful means to develop a hyperpolyhedron within the hyperspace for data mining. The philosophy of the computation methods of "Materials Designer" are as follows:

- (1) It is assumed that the sample points distributed in hyperspace for data mining can be classified into two classes: class "1" and class "2". We usually call class "1" as "good" class if the samples can be classified as "good" or "bad". The distribution zone of sample points of class "1" is called "optimal zone".
- (2) It is assumed that there is only one optimal zone in the hyperspace for data mining. But this assumption is not a limitation of our method, since we have a "LOCAL VIEW" technique which can be used to divide the hyperspace into several subspaces, and we can make each subspace contain only one optimal zone, and make modelling separately by using the boundaries of the subspace as boundary conditions in modelling.
- (3) It is assume that the optimal zone in hyperspace can be described by a hyperpolyhedron with a series of hyperplanes as its boundaries (the hyperpolyhedron can be concave one or convex one. A so-called "box class 2" operation can be used to define a concave hyperpolyhedron in hyperspace). The inequalities describing these hyperplane boundaries can be used as mathematic model for optimal zone.

LMAP is a linear projection method with better separability than traditional pattern recognition methods. It consists of four steps: in the first step the origin of coordinate is moved to the center of the optimal zone. The second step is to generate a hyper-ellipsoid to enclose all sample points of class "1", the third step is to

make the deformation of the ellipsoid into a hypersphere. The final step is to project the figures onto the plane by PCA method.

MREC and ENVELOP methods form hyperpolyhedron in hyperspace to enclose all sample points automatically. By MREC method, various projection maps can be displayed on computer screen to provide useful information of the data structure for you.

In software "Materials Designer", pattern recognition methods are used to reduce the dimensionality of data structure to diminish the number of inputs of ANN and the independent variables of nonlinear regression. In this respect, MREC is most useful tool for dimension reduction.

"MATERIALS DESIGNER" AS A TOOL FOR NEW MATERIALS OPTIMAL DESIGN [1,2,3]

Optimal Design of the Composition of Rare-Earth Containing Phosphor

Keise Optonix Ltd in Germany has applied 45 German patents about the composition of rare-earth-containing phosphors. We have used these data to make data mining by "Materials Designer". The result of data mining indicates that the optimal zone can be extended by extrapolation. By extrapolation we have obtained a series of new compositions located outside of the scope of German patent. Our experimental work confirms that these newly designed phosphores exhibit higher brightness than the German patent declared.

Optimal Design of High Temperature Superconductor

By fluoride dopant addition, the Bi-based superconductor materials has achieved the highest critical temperature to about 116 K. We have used these data to do data mining work by "Materials Designer". Based on the result of data mining, some new composition and sintering condition has been proposed, and the critical temperature level achieved has been elevated to 121 K.

Optimal Design of VPTC ceramic semiconductors

VPTC materials, a special kind of ceramic semiconductor, have been prepared. One of the performance parameter is the ratio between the electric resistance at 0 °C and minimum resistance. The highest ratio achieved is 20. By using "Materials Designer", some proposed new composition and technological condition give much better result: this ratio is elevated to 27.3.

"MATERIALS DESIGNER" AS A TOOL FOR UNKNOWN PHASE DIAGRAM PREDICTION

Computerized Prediction of Ternary Intermetallic Compounds

It is well known that thermodynamic method is widely used to predict ternary phase diagrams based on the data of known binary phase diagrams. But the result is reliable only in the case that there is no unknown ternary intermediate compound formation, since the existence of ternary intermediate compound cannot be predicted by thermodynamic method. The use of "Materials Designer" may solve this problem, because we can use it to find the regularities of the formation of ternary compounds by the data mining of the data of known phase diagrams in the hyperspace spanned by suitable atomic parameters or their functions.

The prediction of ternary intermetallic compounds can be cited as an example. Based on the data from the Data Bank of Ternary Alloy Systems edited by Villars, we select the data of known phase diagrams of 2400 ternary alloy systems as training set, and use Villars's system of atomic parameters (VE , X_{mb} , R_{sp}) or their functions to span hyperspaces. Some very good regularities can be found if the systems consisting of nontransition metals and that consisting of transition metals are treated separately. Based on the regularities found (here the zone of class "1" is the zone of ternary compound forming zone), a series of newly discovered ternary intermetallic compounds can be "predicted" in this way.

Computation of the Liquidus Surfaces of Some Ternary Phase Diagrams of Halide Systems

Based on the data of known ternary phase diagrams of halide systems, we can use "Materials Designer" to make data mining within some hyperspaces spanned by the values of the ionic radii, the ionic charge and the electronegativities of constituent elements. The position of the contour curves of the liquidus surfaces of many ternary systems (for example, the contour curve of the liquidus surface of RbCaCl_3 in CaCl_2 - BaCl_2 - RbCl system) can be predicted with good result).

REFERENCES

1. Chen Nianyi, 1988. Chemical pattern recognition research in China, *Analytica Chimica Acta*, 210(1-2), 175-179.
2. Chen Nianyi, Li Chonghe, Yao Shuwen, Wang Xueye, 1996. Regularities of melting behaviours of some binary alloy phases. *Journal of Alloys and Compounds*, 234, 126-136
3. Chen Nianyi, Li Chonghe, Liu Gang, Qin Pei, 1996. On the formation of ternary alloy phases, *Journal of Alloys and Compounds*, 245, 179-187

AUTHOR'S INDEX

M.F. Abbod	215	N. Chen	1381
S.M. Adballah	1017	N. Chen	1419
S. Ahn	1047	R. Chen	1419
J.M. Abe	695	Y.-M. Chen	1061
A. Akhtar	417	D. Cheung	1079
J. Ahola	531	C.-C. Chiang	1131
N. Aikawa	607	D.J. Choo	947
J.D. Allen, Jr.	961, 989	C.Y. Chung	1023
K. Ameyama	1041	S. Cierpisz	933
W. Andreoni	1397	D.J. Clancy	871
A. Arioti	629	E.J. Colville	649
J.F. Atkinson	347	J.A. Cooper	191
		C.S. Cornelius	325
G. Baiden	53	L. Cser	531
T.J. Bailey	921	C. Curtis	317
J. Balcita	499		
P. Barr	111	J. Daams	1339
R. Barton	269	M. d'Amore	703
O.A. Bascur	829	D. Dasgupta	257
D. Bassi	975	W.J. Davis	615
M. Benedict	1185	L.R.P. De Andrade Lima	505
R.R. Biggers	1258, 1317	D.V. Dempsey	1258, 1317
Y. Bissiri	635	S. Dessureault	145
H. Bode	339	R.J. Dippenaar	75
G. Bonifazi	465, 485	S. Dolinšek	847
B.M. Brasfield	347	A. Donnarumma	185, 663
J.C. Bressiani	797	R. Doraiswami	735
R.T. Bui	749	B.F. Duan	1361
J.D. Busbee	1258, 1317	M. Duarte	975
		K. Dudek	543
T.L. Calton	347	S. Dunbar	145, 635
J.J. Campbell	939	M.N. Durakbasa	927
N. Cappetti	185		
M.J. Cardew-Hall	1017	S.A. Ehikioya	139
L.-E. Carlsson	459	H. Eldeib	447
J.C. Cassa	291, 381	J. Endou	817
O. Castillo	151, 855	J.R. Esslinger	331
A.C.D. Chaklader	797	R.N. Evans	331
T. Chandra	105		
T. Chashikawa	453	M. Fabiunke	655
C.C. Chang	789	C. Fantozzi	629
J.Y. Chen	901	N. Farmer	879
L.S. Chen	805	P. Farrington	157
M.Y. Chen	395	M. Fathi-Torbaghan	1011

F. Ferguson	317	S. Hirose	233
R. Felix	299	B. Hlaváček	403
M. Ferry	105	C.T.T. Ho	1061
G. Floridia	291, 381	P.D. Hodgson	389, 953
G.A. Fodor	895	D.A. Holder	157
S. Forouzi	967	D.A. Holder	347
S. Forrest	257	R.-Q. Hsu	1029
J.M. Fragomeni	577, 585	C.-C. Hu	1131
W.G. Frazier	1139	H.M. Huang	285
K. Fujii	453	I.B. Huang	423
S. Fuks	1123	W. Huang	1277
M. Furukawa	1115	P. Hubík	437
Y. Furukawa	21	G. Huh	947
Y. Fukuhara	743	Y.S. Hwang	879
		H. Hyötyniemi	11, 179, 459
M. Geiger	641		
L.M. Geng	91	B. Igelnik	367
R. Gerth	1151	K. Ishida	373
D.T. Gethin	513, 1035	K. Ishino	1093
M.M. Ghomshei	519	N. Ivezic	961, 989
H. Ghulman	1151	S. Iwata	1323, 1399
D.A. Gibson	325		
P. Giordano	703	A.G. Jackson	1185
Z. Gomolka	813	J. Jang	947
G.D. González	59	H-G. Jeong	429
M. Granchi	629	P.D. Jero	1241
J.L. Grantner	895	L. Jin	805
R.W. Grimes	1197	J.G. Jones	1241, 1258, 1317
W.A. Gruver	839	H.K. Jung	593
A. Grzech	823	K.D. Jung	593
C. Guist	339		
S.R. Gunn	361	B. Kádár	131
M.J. Guo	779	R. Kainuma	373
Y.M. Guo	861	K. Kamitani	565
M. Gupta	119	J.S. Kandola	361
		C.G. Kang	593
A. Hambaba	1073	S. Kang	1047
S. Hanada	429	A. Karcher	623
R.D. Harrell	157, 347	A. Katayama	607
J. Hart	879	M. Kato	373
Y. Hasegawa	695	K. Katoh	571
J. Hätönen	459	Y. Kawazoe	..355
H. Helman	537, 549, 561	A.R. Khoei	513, 1035
L. Hildebrand	1011	H.S. Kim	429
T. Hirasawa	221	J. Kim	1263
T. Hirasawa	245	N.N. Kiselyova	1387

M. Kitabata	995	P. Ma	867
Gerd Kock	655	S. Mackinson	491
T. Kohonen	27	J.F. Maguire	1235
N. Koga	403	J.T. Malin	1179
L.X. Kong	389, 953	K. Manabe	571, 601, 607
J. Kopac	847	P.A. Manohar	105
A.S. Korhonen	531	A. Mansour	1103
H. Kotera	221, 245, 565	P.Mäntylä	531
H.Koyama	571	J.J.Mareš	437
G. Kozłowski	1258, 1317	P. Massacci	485
J. Krištofik	437	I. Masters	513
C. Kropas-Hughes	1305	E. Medina	1157
C.C. Kung	681	J.A. Meech	111, 309, 499, 519, 975
J.A. Kurien	871	K.J. Meech	445
A. Kusiak	887	G. Meghabghab	729
J. Kusiak	543, 773	B. Mehta	1151
D-W. Kum	429	P. Melin	151, 855
T. Kurita	239	T. Menzel	641
K. Kyuma	1297	S. Messimer	157
		J. Metcalfe	215
		J. Miettunen	459
W.C. Lai	681	L. Miller-Tait	983
K.C. Lau	1023	L. Monostori	131, 847
S.R. LeClair	367, 1235, 1263, 1361, 1399	C. Mui	111
E.S. Lee	1047	H. Munekata	233
J.Y. Lee	947		
P.D. Lee	1197	Y. Nagasaka	85
R.S. Lee	1061	S. Nagatomo	565
Y.J. Lee	1285	K. Nakamatsu	695
Y.K. Lee	947	K. Nakanishi	861
J.G. Lenard	543	G. Nasr	729
J. Leopold	227	G.M. Nicoletti	713, 1087
P.L. Leung	285	A.J. Niemi	11
R.W. Lewis	1035	H. Nishimura	601
J.Y. Li	765	D.H. Norrie	887
K.P. Li	91	B. Novak	201
Y.P. Li	867	B. Nichols	215
G. Lin	687, 1003		
L. Lin	953	H. Ohtani	373
S.-C. Lin	669	S. Oka	1277
Y. Liu	939	R.T. Oliveira	291, 381
D.A. Linkens	215, 395, 755, 921	N. Ono	373
H.R. Liu	765	T. Ono	239
R.S. Liu	765	C. Orłowski	195
W. Lu	1419	P.H. Osanna	927
Y. Luo	867	M. Oxley	1371

E.P. Paladini	165	B. Shi	687, 1003
R. Pakalnis	983	Y. Shigaki	561
Y.-H. Pao	367, 1361	S. Shima	221, 245, 563
M. Pappalardo	185, 663, 703	F. Shimaya	555
Z. Pawlak	37	M. Shinkawa	555
G.H. Park	1285	S.S. Shivathaya	105
G.R. Park	1263	G.R. Shumaker	1163
J. H. Park	1285	O. Simula	531
S.H. Park	947	I. Sinclair	361
P.E. Parker	895	A.K. Sirén	477
C.D.M. Pataro	549	M.H. Smith	839
A. Pellegrino	703	G.W. Snyder	331
J. Perron	749	S.J. Spencer	939
M. Pietrzyk	773	B. Štěpánek	437
J. Poindexter	1163	O. Stephansson	471
		M. Stevenson	735
P. Qin	1419	A.R. Souza	291, 381
M.V. Quintella-Cury	1123	Z. Strnad	431
		B.A. Stucke	1163
C. Reidsema	1055	A. Suárez	975
S. Reimann	1103	Y. Suehiro	1041
S. Rajan	735	S. Sugiyama	909
P.A.S. Reed	361	A. Suzuki	695
D.A. Ress	1225	H. Szczerbicka	269
D. Rochowiak	157	3 E. Szczerbicki	813, 1055
J. Rogers	157		
B.F. Rolfe	1017	M. Tabib-Azar	367, 1249
J.A. Romagnoli	947	Y. Takefuji	453, 723, 743, 995,
D. Roy	281	1109, 1271, 1277	
P. Russo	703	A. Talaie	947
D. Russell	157	J. Tenner	921
		L. Tikasz	749
Y. Sahai	119	K.M. Tiwari	281
T. Saito	723	S.K. Tiwari	281
Y. Saito	555	A. Torres	309
Y. Sakamoto	221, 245	V. Torres	309
G. Samadi Hosseinali	417	T.T. Tran	139
I.V. Samarasekera	73	D. Tromans	411
S. Sändig	251	E.G. Truelove	139
E. Santoro	185	U.-B. Tsai	1029
M. Sato-Ilic	207	H.-L. Tsoi	1079
M. Scoble	145	I.B. Turksen	173
H.-J. Sebastian	..163		
S. Serranti	465	O. Unold	277
J. Šesták	403, 431, 437	T. Ushio	299
V. Šestáková	437	H. Utsunomiya	555

T. Van Le	675	T. Yamamoto	221
M.M. Veiga	797	M. Yamaura	233
S.M.B.Veiga	797	M. Yang	607
E. Vettori	629	Y.Y. Yang	755
N.K. Vidyarthi	281	S.K. Yen	423, 779, 789
Z.J. Viharos	847	R. Ylinen	11
P. Villars	1339	A. Yoshida	221
P. Villars	1399	S. Yoshihara	571, 601
R.C. Villas Bôas	481, 505	N. Yoshiike	1109
F. Volpe	465	G.T. Yu	785
		W.S. Yu	1125
R.H. Wagoner	91		
B. Wang	953	M.O. Zacate	1197
D. Wang	805	L.A. Zadeh	3
W.X. Wang	471	T. Zacharia	961, 989
Z. Wang	1067	H.Z. Zan	779
M. Watanabe	1115	L.E. Zárate	537
K.R. Weller	939	B.L.Zhang	867
G.A.W. West	1017	T. Zhang	839
P. Wiesner	251	X. Zhang	887
M. Williams	785	Y.L. Zhao	1361
J. Wirtz	623	B.S. Zhu	867
C.-W. Wu	1029	D.D. Zhu	1381
L Xie	1067	H.J. Zimmermann	45
Y. Xu	805	R. Zuco	465